



Munich Personal RePEc Archive

# **Coevolution of Deception and Preferences: Darwin and Nash Meet Machiavelli**

Yuval Heller and Erik Mohlin

Univeristy of Oxford

13. August 2015

Online at <https://mpra.ub.uni-muenchen.de/66177/>

MPRA Paper No. 66177, posted 23. August 2015 13:35 UTC

# Coevolution of Deception and Preferences: Darwin and Nash Meet Machiavelli\*

Yuval Heller<sup>†</sup> and Erik Mohlin<sup>‡</sup>

3/8/15

## Abstract

We develop a framework in which individuals' preferences coevolve with their abilities to deceive others about their preferences and intentions. Specifically, individuals are characterized by (i) a level of cognitive sophistication and (ii) a subjective utility function. Increased cognition is costly, but higher-level individuals have the advantage of being able to completely deceive lower-level opponents about their preferences and intentions. Only individuals who are of the same cognitive level can observe each other's preferences. Our main result shows that, despite the limited possibility to observe preferences, and despite the strong form of deception, essentially only efficient population states can be stable. Moreover, if the marginal cognitive costs are not too high, then only efficient Nash equilibria are stable. We extend our model to study preferences that depend also on the opponent's type.

**Keywords:** Evolution of Preferences; Indirect Evolutionary Approach, Theory of Mind; Depth of Reasoning; Deception; Efficiency. **JEL codes:** C72, C73, D03, D83.

## 1 Introduction

For a long time economists took preferences as given. The study of their origin and formation was considered a question outside the scope of economics. Over the past two decades this has changed dramatically. In particular, there is now a large literature on the evolutionary foundations of preferences (for an overview, see Robson and Samuelson, 2011). A prominent strand of this literature is the so-called “indirect evolutionary approach”, pioneered by Güth and Yaari (1992) (term coined by Güth, 1995). This approach has been used to explain the existence of a variety of “non-standard” preferences that do not coincide with material payoffs, e.g., altruism, spite, and reciprocal preferences.<sup>1</sup> Typically, the non-materialistic preferences in question convey

---

\*Valuable comments were provided by Vince Crawford, Itzhak Gilboa, Larry Samuelson, and Jörgen Weibull, as well as participants at presentations in Oxford, at G.I.R.L.13 in Lund, the Toulouse Economics and Biology Workshop, DGL13 in Stockholm, the 25th International Conference on Game Theory at Stony Brook, and the Biological Basis of Preference and Strategic Behaviour conference at Simon Fraser University 2015. Erik Mohlin was supported in part by the European Research Council, grant no. 230251.

<sup>†</sup>Affiliation: Queen's College and Department of Economics, University of Oxford. Address: Queen's College, High Street, Oxford, OX1 1AW, United Kingdom. E-mail: yuval.heller@queens.ox.ac.uk.

<sup>‡</sup>Affiliation: Nuffield College and Department of Economics, University of Oxford. Address: Nuffield College, New Road, Oxford OX1 1NF, United Kingdom. E-mail: erik.mohlin@nuffield.ox.ac.uk.

<sup>1</sup>For example, Bester and Güth (1998), Bolle (2000), and Possajennikov (2000) study combinations of altruism, spite, and selfishness. Ellingsen (1997) finds that preferences that induce aggressive bargaining can survive in a Nash demand game. Fershtman and Weiss (1998) study evolution of concerns for social status. Sethi and Somanthan (2001) study the evolution of reciprocity in the form of preferences that are conditional on the opponent's preference type. In the context of the finitely repeated Prisoner's Dilemma, Guttman (2003) explores the stability of conditional cooperation. Dufwenberg and Güth (1999) study firm's preferences for large sales. Güth and Napel (2006) study preference evolution when players use the same preferences in both ultimatum and

some form of commitment advantage that induces opponents to behave in a way that benefits individuals with non-materialistic preferences, as described by Schelling (1960) and Frank (1987). Indeed, Heifetz, Shannon, and Spiegel (2007) show that this kind of result is generic.

A crucial feature of the indirect evolutionary approach is that preferences are explicitly or implicitly assumed to be at least partially observable.<sup>2</sup> Consequently the results are vulnerable to the existence of mimics who signal that they have, say, a preference for cooperation, but actually defect on cooperators, thereby earning the benefits of having the non-standard preference without having to pay the cost (Samuelson, 2001). The effect of varying the degree to which preferences can be observed has been investigated by Ok and Vega-Redondo (2001), Ely and Yilankaya (2001), Dekel, Ely, and Yilankaya (2007), and Herold and Kuzmics (2009). They confirm that the degree to which preferences are observed decisively influences the outcome of preference evolution.

However, the degree to which preferences are observed is still exogenous in these models. In reality we would expect both the preferences, and the ability to observe or conceal them, to be the product of an evolutionary process.<sup>3</sup> *This paper studies the missing link between evolution of preferences and evolution of how preferences are concealed, feigned, and detected.* In our model the ability to observe preferences and the ability to deceive and induce false beliefs about preferences are endogenously determined by evolution, jointly with the evolution of preferences. Cognitively more sophisticated players completely deceive cognitively less sophisticated players. Mutual observation of preferences occurs only when players of the same cognitive level meet. *Despite allowing for only a “grain” of mutual observability, and despite assuming a strong form of deception, we find that, generically, only efficient outcomes can be played in stable population states.* Moreover, if the marginal cognitive cost is not too high, then only Nash outcomes can be played in stable population states. Conversely, we find that, generically any efficient Nash equilibrium can be implemented in a stable population state. Thus, *essentially, an outcome is stable if and only if it is an efficient Nash equilibrium.*

**Overview of the Model.** As in standard evolutionary game theory we assume an infinite population of individuals who are uniformly randomly matched to play a symmetric normal form game.<sup>4</sup> Each individual has a type, which is a tuple, consisting of a *preference component* and a *cognitive component*. The preference component is identified with a utility function over the set of outcomes (i.e., action profiles). In an extension we allow for *type-interdependent preferences*, which are represented by utility functions that are defined over both action profiles and the opponent’s type. The cognitive component is simply a natural number, representing the

---

dictator games. Koçkesen and Ok (2000) investigate survival of more general interdependent preferences in aggregative games. Friedman and Singh (2009) show that vengefulness may survive if observation has some degree of informativeness. Recently, Norman (2012) show how to adapt some of these results into a dynamic model

<sup>2</sup>Gamba (2013) is an interesting exception. She assumes play of a self-confirming equilibrium, rather than a Nash equilibrium, in an extensive form game. This allows for evolution of non-materialistic preferences even when they are completely unobservable. An alternative is to allow for a dynamic that is not strictly payoff monotonic. This approach is pursued by Frenkel, Heller, and Teper (2014), who show that multiple biases (inducing non-materialistic preferences) can survive in non-monotonic evolutionary dynamics even if they are unobservable, because each approximately compensates for the errors of the others.

<sup>3</sup>On this topic, Robson and Samuelson (2011) write: “The standard argument is that we can observe preferences because people give signals – a tightening of the lips or flash of the eyes – that provide clues as to their feelings. However, the emission of such signals and their correlation with the attendant emotions are themselves the product of evolution. [...] We cannot simply assume that mimicry is impossible, as we have ample evidence of mimicry from the animal world, as well as experience with humans who make their way by misleading others as to their feelings, intentions and preferences. [...] In our view, *the indirect evolutionary approach will remain incomplete until the evolution of preferences, the evolution of signals about preferences, and the evolution of reactions to these signals, are all analysed within the model.*” [Emphasis added] (pp. 14–15)

<sup>4</sup>It is known that positive assortative matching is conducive to the evolution of altruistic behaviour (Hines and Maynard Smith, 1979) and non-materialistic preferences even when preferences are perfectly unobservable (Alger and Weibull, 2013, Bergstrom 1995). It is also known that finite populations allow for the evolution of spiteful behaviours (Schaffer, 1988) and non-materialistic preferences (Huck and Oechssler, 1999). By assuming that individuals are uniformly randomly matched in an infinite population, we avoid confounding these effects with the effect of endogenising the degree of observability.

level of cognitive sophistication of the individual.<sup>5</sup> The cost of increased cognition is strictly positive.

When the individuals in a match are of *different* cognitive levels, the one with the higher level is assumed to be able to deceive the one with the lower level. For the sake of tractability, and in order to “stack the cards” against our main result, we model a strong form of deception. The deceiver observes the opponent’s preferences perfectly, and is allowed to choose whatever she wants the deceived party to believe about the deceiver’s intended action choice. A strategy profile that is consistent with this form of deception is called a *deception equilibrium*. When both individuals are of the *same* cognitive level, we assume that each player observes the opponent’s preferences, and, as a result, the individuals play a Nash equilibrium of the complete information game induced by their preferences.<sup>6</sup>

The state of a population is described by a *configuration*, consisting of a type distribution and a behaviour policy. The *type distribution* is simply a finite support distribution on the set of types. The *behaviour policy* specifies a Nash equilibrium for each match between cognitive equals, and a deception equilibrium for each match between types of different cognitive levels. In a *neutrally stable configuration* all incumbents earn the same, and if a small group of mutants enter they earn weakly less than the incumbents in any *focal* post-entry state. A focal post-entry state is one in which the incumbents behave against each other in the same way as before the mutants entered.

**Main Results.** We say that a profile is efficient if it maximizes the sum of fitness payoffs. Theorem 1 demonstrates that in any stable configuration, any type  $\bar{\theta}$  with the highest cognitive level in the incumbent population must play efficiently when meeting itself. The intuition is that otherwise a highest-type mutant who mimics the play of  $\bar{\theta}$  against all incumbents while playing efficiently against itself, would outperform type  $\bar{\theta}$  (an application of the “secret handshake” argument due to Robson, 1990).

Next we restrict attention to generic games (i.e., games in which each fitness payoff is independently drawn from a continuous distribution) and obtain our main result: *any stable configuration must induce efficient play* in all matches between all types. The idea of the proof can be briefly sketched as follows. We first show that any type  $\theta$  in a stable configuration must play efficiently when meeting *itself*. Otherwise a mutant who has the same level as  $\theta$  and the same utility function as  $\theta$ , but who plays efficiently against itself, could invade the population. Next, we show that *any* two types must play efficiently. The intuition is that otherwise the mean within-group fitness would be higher than the between-group fitness, which implies that there would be instability to small perturbations in the frequency of the types: A type who became slightly more frequent would have a higher fitness than the other incumbents, and this would take the population away from the original configuration.

The existing literature (e.g., Dekel, Ely, and Yilankaya, 2007) has demonstrated that if players perfectly observe the opponent’s preferences (or do so with sufficiently high probability), then only efficient outcomes are stable. *Our key contribution is to show that a “grain” of perfect observability is enough to ensure that stability implies efficiency.* Unlike the existing models, which assume that any player perfectly observes the preferences of any other player, our model assumes perfect observability only as a “tie-breaking” rule: players with equal levels perfectly observe each other’s preferences, but a player who acquires a higher cognitive level (which may incur an arbitrarily low cognitive cost) is able to completely deceive her opponent about her preferences.

Next we make a few observations, which allow us to fully characterize stable configurations in generic games.

---

<sup>5</sup>The one-dimensional representation of cognitive ability reflects the idea that if one is good at deceiving others, then one is more likely to be good also at reading others and avoiding being deceived by them. In this paper we simplify this relation by assuming a perfect correlation between the two abilities, and leave the study of more general relations for future research.

<sup>6</sup>Our assumption that players of equal levels play a Nash equilibrium can be motivated by the (so-called) “folk theorem of evolutionary game theory”, which implies that, for many evolutionary dynamics, any stable population state is a Nash equilibrium of the underlying complete information game (see Nachbar, 1990, for a formal statement and proof).

We note that any generic game admits at most one efficient symmetric action profile, and that this symmetric action profile is either a strict equilibrium or not an equilibrium at all. In the former case, we demonstrate the stability of a homogeneous population of individuals with the lowest cognitive level and with preferences such that the efficient action is dominant. Moreover, if the cognitive cost required for achieving the second cognitive level is sufficiently low, then this stable configuration is essentially unique. In any other case, i.e., where the underlying game does not have an efficient symmetric action profile, or if an efficient symmetric action profile exists but is not an equilibrium, then there is no stable configuration. We also show that the same characterization holds for pure stable configurations in non-generic games.

Finally, we note that non-generic games may admit different kinds of stable configurations. One particularly interesting family of non-generic games is the family of zero-sum games, such as the Rock-Paper-Scissors game. We analyse this game and characterize a heterogeneous stable population (inspired by a related construction in Conlisk, 2001) in which different cognitive levels coexist, players with equal levels play the Nash equilibrium, and players with higher levels beat their opponents but this gain is offset by higher cognitive costs.

**Variants and Extensions.** As mentioned above, our main result relies on having perfect observability as the “tie-breaking” rule in matches between cognitive equals. In Section 5 we consider the opposite assumption, namely, that players with equal levels do not observe each other’s preferences. This means that there are no matches in which there is mutual observation of preferences. We show that this variant yields very different results. Specifically, we show that (1) any pure strict (possibly inefficient) equilibrium of the underlying game is stable, and, as a partial converse, (2) the individuals of the highest type always play a Nash equilibrium (again, not necessarily efficient) when being matched against themselves. Thus, whether stability implies efficiency or Nash equilibrium behaviour, crucially depends on the assumption of a “grain” of exogenous observability in the matches between equals in our model. The existing literature (e.g., Ok and Vega-Redondo, 2001; Ely and Yilankaya, 2001; Dekel, Ely, and Yilankaya, 2007; Norman, 2012) shows that if players cannot observe each other’s preferences, then only Nash equilibria can be stable. *Our results show that non-observability among players with equal levels is enough to imply that a strict Nash equilibrium is sufficient for stability*, even though our setup allows a player to spend a (possibly arbitrarily low) additional cognitive cost, in order to obtain a higher level than her opponent, and perfectly observe her opponent’s preferences.

In most of the paper we deal only with “type-neutral” preferences that are defined only over action profiles. Section 6 extends the analysis to interdependent preferences, i.e., preferences that may also depend on the opponent’s type. Herold and Kuzmics (2009) study a similar setup while assuming perfect observability of types among all individuals. Their key result is that any mixed action that gives each player a payoff above her maxmin payoff can be the outcome of a stable configuration.<sup>7</sup> We obtain a slightly weaker extension of this result while assuming only perfect observability among individuals with equal cognitive levels. Specifically, we show that stable pure population states (i.e., populations in which everyone plays the same pure action) are essentially Nash equilibria that yield each player a payoff above the maxmin value. We conclude by characterizing stable configurations in the “Hawk-Dove” game (Section 6.4). We show that such games admit heterogeneous stable configurations in which players with different levels coexist, each type has discriminating preferences that induce

---

<sup>7</sup>Herold and Kuzmics (2009) expand the framework of Dekel, Ely, and Yilankaya (2007) to include interdependent preferences, i.e., preferences that depend on the opponent’s preference type. Under perfect or almost perfect observability, if all preferences that depend on the opponent’s type are considered, then any symmetric outcome above the minmax material payoff is evolutionarily stable. In our setting a pure profile also has to be a Nash equilibrium in order to be the sole outcome supported by evolutionarily stable preferences. Herold and Kuzmics (2009) find that non-discriminating preferences (including selfish materialistic preferences) are typically not evolutionarily stable on their own. By contrast, certain preferences that exhibit discrimination are evolutionarily stable. Similarly, evolutionary stability requires the presence of discriminating preferences also in our setup.

cooperation only against itself, and higher types “exploit” lower types (and this is offset by their higher cognitive levels).

In Appendix B we briefly present another variant of the model in which the deceiver is *not* able to “tailor” the attempted deception to the current opponent’s type. Instead, an individual has to use the same attempted deception against all opponents. It turns out that qualitatively similar results hold with this less flexible form of deception.

**Further related literature.** There is a large literature in biology and evolutionary psychology on the evolution of “theory of mind” (Premack and Woodruff 1979). According to the “Machiavellian intelligence” hypothesis (Humphrey, 1976), and “social brain” hypothesis (Dunbar, 1998), the extraordinary cognitive abilities of humans evolved as a result of the demands of social interactions, rather than the demands of the natural environment: in a single-person decision problem there is a fixed benefit of being smart, but in a strategic situation it may be important to be smarter than the opponent. From an evolutionary perspective, the potential advantage of a better theory of mind has to be traded off against the cost of increased reasoning capacity. Increased cognitive sophistication in the form of higher-order beliefs is associated with non-negligible costs (Holloway, 1996, Kinderman, Dunbar, and Bental 1998). Our model incorporates these features.

There is a smaller literature on the evolution of strategic sophistication within game theory; see, e.g., Stahl (1993), Banerjee and Weibull (1995), Stennek (2000), Conlisk (2001), Abreu and Sethi (2003), Mohlin (2012), Rtischev (2012), and Heller (2015). As in these papers, we provide results to the effect that different degrees of cognitive sophistication may coexist.

Kimborough, Robalino, and Robson (2014) construct a model to demonstrate the advantage of having a theory of mind (understood as an ability to ascribe stable preferences to other players) over learning by reinforcement. In novel games the ascribed preferences allow the agents with a theory of mind to draw on past experience whereas a reinforcement learner without such a model has to start over again. Hopkins (2014) explains why costly signaling of altruism may be especially valuable for those agents who have a theory of mind.

Robson (1990) initiated a literature on evolution in cheap talk games by formulating the secret handshake effect: evolution selects an efficient stable state if mutants can send messages that the incumbents either do not see or not benefit from seeing. Against the incumbents a mutant plays the same action as the incumbents do, but against other mutants the mutant plays an action that is a component of the efficient equilibrium. Thus the mutants are able to invade unless the incumbents are already playing efficiently. See also Matsui (1991). As pointed out by Wärneryd (1991), and Schlag (1993), among others, problems arise if either the incumbents use all available messages (so that there is no message left for the incumbents to coordinate on) or the incumbents follow a strategy that induces the mutants to play an action that lowers the mutants’ payoffs below those of the incumbents. Kim and Sobel (1995) use stochastic stability arguments, and Wärneryd (1998) uses complexity costs, to circumvent this problem. Similarly, evolution selects efficient outcome in our model, where the preferences also serve the function of messages

**Structure.** The rest of the paper is organized as follows. Section 2 presents the model. In Section 3 we define our stability notions. The results for the main model are presented in Section 4. Section 5 deals with a variant in which a player cannot observe the preferences of an opponent with the same cognitive level. Section 6 extends the model to include type-interdependent preferences. Section 7 concludes. Appendix A contains proofs not in the main text. Appendix B presents a variant of the model with uniform deception. Appendix C formally constructs heterogeneous stable populations in Rock-Paper-Scissors and Hawk-Dove games.

## 2 Model

We consider a large population of agents, each of which is endowed with a type that determines her subjective preferences and her cognitive level. The agents are randomly matched to play a symmetric two-player game. A dynamic evolutionary process of cultural learning, or biological inheritance, increases the frequency of more successful types. In the next section, we present a static solution concept to capture stable population states in such environments.

### 2.1 Underlying Game and Types

Consider a symmetric two-player normal form game  $G$  with a finite set  $A$  of pure actions and a set  $\Delta(A)$  of mixed actions (or strategies). We use the letter  $a$  ( $\sigma$ ) to describe a typical pure action (mixed action). Payoffs are given by  $\pi : A \times A \rightarrow \mathbb{R}$ , where  $\pi(a, a')$  is the payoff to a player using action  $a$  against action  $a'$ . The payoff function is extended to mixed actions in the standard way, where  $\pi(\sigma, \sigma')$  denotes the material payoff to a player using strategy  $\sigma$ , against an opponent using strategy  $\sigma'$ . With a slight abuse of notation let  $a$  denote the degenerate mixed strategy that puts all weight on pure strategy  $a$ . We adopt this convention for probability distributions throughout the paper.

*Remark 1.* The restriction to symmetric games is without loss of generality when dealing with interactions in a single population. In cases in which the interaction is asymmetric, it can be captured in our setup (as is standard in the literature; see, e.g., Selten, 1980 and Samuelson, 1991) by embedding the asymmetric interaction in a larger, symmetric game in which nature first randomly assigns the players to roles in the asymmetric interaction.

We imagine a large (technically infinite) population of individuals who are uniformly randomly matched to play the game  $G$ . Each individual  $i$  in the population is endowed with a *type*  $\theta = (u, n) \in \Theta = U \times \mathbb{N}$ , consisting of *preferences*, identified with a von Neumann-Morgenstern utility function,  $u \in U$  and a *cognitive level*  $n \in \mathbb{N}$ . Let  $\Delta(\Theta)$  be the set of all finite support probability distributions on  $\Theta$ . A population is represented by a finite support *type distribution*  $\mu \in \Delta(\Theta)$ . Let  $C(\mu)$  denote the support (carrier) of type distribution  $\mu \in \Delta(\Theta)$ . Elements of  $C(\mu)$  will be called incumbents. Given a type  $\theta$ , we use  $u_\theta$  and  $n_\theta$  to refer to its preferences and cognitive level, respectively.

In the main model we assume that the preferences are defined over action profiles, as in Dekel, Ely, and Yilankaya (2007).<sup>8</sup> This means that any preferences can be represented by a utility function of the form  $u : A \times A \rightarrow \mathbb{R}$ . The set of all possible (modulo affine transformations) utility functions on  $A \times A$  is  $U = [0, 1]^{|A|^2}$ . Let  $BR_u(\sigma')$  denote the set of best replies to strategy  $\sigma'$  given preferences  $u$ , i.e.,  $BR_u(\sigma') = \arg \max_{\sigma \in \Delta(A)} u(\sigma, \sigma')$ .

There is a fitness cost to increased cognition, represented by a strictly increasing cognitive cost function  $k : \mathbb{N} \rightarrow \mathbb{R}_+$ . The fitness payoff of an individual equals the material payoff from the game, minus the cognitive cost. Let  $k_n$  denote the cost of having cognitive level  $n$ . Hence  $k_\theta = k_{n_\theta}$  denotes the cost of having type  $\theta$ . Without loss of generality, we assume that  $k_1 = 0$ . In some of our results we will make the additional assumption that  $k_2$  is sufficiently small.

### 2.2 Configurations

A state of the population is described by a type distribution and a behaviour policy for each type in the support of the type distribution. An individual's behaviour is assumed to be (subjectively) rational in the sense

<sup>8</sup>In Section 6, we study *type-interdependent* preferences, which depend also on the opponent's type, as in Herold and Kuzmics (2009).

that it maximizes her subjective preferences given the belief she has about the opponent’s expected behaviour. However, her beliefs may be incorrect, if she is deceived by her opponent. An individual is deceived if and only if her opponent is of a higher cognitive level.

If two individuals of the same cognitive level are matched to play, then they play a Nash equilibrium of the game induced by their preferences. Given two preferences  $u, u' \in U$ , let  $NE(u, u') \subseteq \Delta(A) \times \Delta(A)$  be the set of mixed equilibria of the game induced by the preferences  $u$  and  $u'$ , i.e.,

$$NE(u, u') = \{(\sigma, \sigma') \in \Delta(A) \times \Delta(A) : \sigma \in BR_u(\sigma') \text{ and } \sigma' \in BR_{u'}(\sigma)\}.$$

*Remark 2.* Similar to most on the existing literature of the indirect evolutionary approach (e.g., Güth and Yaari, 1992; Dekel, Ely, and Yilankaya, 2007, Section 3), we assume perfect observability of the opponent’s preferences. However, while the existing literature assumes this for all interactions, we assume it only when two incumbents with the same cognitive level interact. When two agents with different levels meet, the observability is endogenously determined by the deception efforts of the player with the higher level, and the observability assumption is only used as a “tie-breaking rule”. In Section 5 we analyse the opposite “tie-breaking rule”, according to which, when agents with the same cognitive level interact, they do not observe the opponent’s preferences.

If two individuals of different cognitive levels are matched to play, then the individual with the higher cognitive level (henceforth, the *higher type*) observes the opponent’s preferences perfectly, and is able to deceive the opponent (henceforth, the *lower type*). The deceiver is allowed to choose whatever she wants the deceived party to believe about the deceiver’s intended action choice. The deceived party best-responds given her possibly incorrect belief.

For simplicity, we assume that if the deceived party has multiple best replies, then the deceiver is allowed to break indifference, and choose which of the best replies she wants the deceived party to play. Consequently the deceiver is able to induce the deceived party to play any strategy that is a best reply to some belief about the opponent’s mixed action, given the deceived party’s preferences.

Given preferences  $u \in U$ , let  $\Sigma(u)$  denote the set of *undominated strategies*. By the minmax theorem,  $\Sigma(u)$  is also the set of actions that are best replies to at least one strategy of the opponent (given the preferences  $u$ ). Formally, we define

$$\Sigma(u) = \{\sigma \in \Delta(A) : \text{there exists } \sigma' \in \Delta(A) \text{ such that } \sigma \in BR_u(\sigma')\}.$$

We say that a strategy profile is a *deception equilibrium* if the strategy profile is optimal from the point of view of player  $i$  under the constraint that player  $j$  has to play an undominated strategy. Formally:

**Definition 1.** Given two types  $\theta, \theta'$  with  $n_\theta > n_{\theta'}$ , a strategy profile  $(\tilde{\sigma}, \tilde{\sigma}')$  is a *deception equilibrium* if

$$(\tilde{\sigma}, \tilde{\sigma}') \in \arg \max_{\sigma \in \Delta(A), \sigma' \in \Sigma(u_{\theta'})} u_\theta(\sigma, \sigma').$$

Let  $DE(\theta, \theta')$  be the set of all such deception equilibria.

We are now in a position to define our key notion of a configuration, by combining a type distribution with a behaviour policy, as represented by Nash equilibria and deception equilibria.

**Definition 2.** A *configuration* is a pair  $(\mu, b)$  where  $\mu \in \Delta(U)$  is a type distribution, and  $b : C(\mu) \times C(\mu) \rightarrow$



$\Delta(A)$  is a *behaviour policy* such that for each  $\theta, \theta' \in C(\mu)$  :

$$n_\theta = n_{\theta'} \implies (b_\theta(\theta'), b_{\theta'}(\theta)) \in NE(\theta, \theta'), \text{ and}$$

$$n_\theta > n_{\theta'} \implies (b_\theta(\theta'), b_{\theta'}(\theta)) \in DE(\theta, \theta').$$

We interpret  $b_\theta(\theta') = b(\theta, \theta')$  as the strategy of type  $\theta$  when being matched with type  $\theta'$ .

Given a configuration  $(\mu, b)$  we call the types in its support the *incumbents*.

Note that standard arguments imply that for any type distribution  $\mu$  there exists a mapping  $b : C(\mu) \times C(\mu) \rightarrow \Delta(A)$  such that  $(\mu, b)$  is a configuration.

The expected fitness to an individual of type  $\theta$  in configuration  $(\mu, b)$  is:

$$\Pi_\theta((\mu, b)) = \sum_{\theta' \in C(\mu)} \mu(\theta') \cdot \pi(b_\theta(\theta'), b_{\theta'}(\theta)) - k_\theta.$$

When all incumbent types have the same expected fitness, we say that the configuration is *balanced*, and denote this uniform expected payoff by  $\Pi((\mu, b))$ .

*Remark 3.* Our model assumes that a player may use different deceptions against different types with lower cognitive levels. We note that all our results remain the same (with minor changes to the proofs) in an alternative setup in which individuals have to use the *same* mixed action in their deception efforts towards all opponents with lower cognitive levels. We refer to this as *uniform deception*. The formal changes in the model that are required to implement this variant are described in Appendix A.

### 3 Evolutionary Stability

Recall that a neutrally stable strategy (Maynard Smith and Price, 1973 and Maynard Smith, 1982) is a strategy that, if played by most of the population, weakly outperforms any other strategy. Similarly, an evolutionarily stable strategy is a strategy that, if played by most of the population, strictly outperforms any other strategy.

**Definition 3.** A strategy  $\sigma \in \Delta(A)$  is a *neutrally stable strategy (NSS)* if for every  $\sigma' \in \Delta(A)$  there is some  $\bar{\varepsilon} \in (0, 1)$  such that if  $\varepsilon \in (0, \bar{\varepsilon})$ , then  $\tilde{\pi}(\sigma', (1 - \varepsilon)\sigma + \varepsilon\sigma') \leq \tilde{\pi}(\sigma, (1 - \varepsilon)\sigma + \varepsilon\sigma')$ . If the weak inequality is replaced by strict inequality for each  $\sigma' \neq \sigma$ , then  $\sigma$  is an *evolutionarily stable strategy (ESS)*.

We extend the notions of neutral and evolutionary stability, from strategies to configurations. We begin by defining the type game that is induced by a configuration.

**Definition 4.** For any configuration  $(\mu, b)$  the corresponding *type game*  $\Gamma_{(\mu, b)}$  is the symmetric two-player game where each player's strategy space is  $C(\mu)$ , and the payoff to strategy  $\theta$ , against  $\theta'$ , is  $\pi(b_\theta(\theta'), b_{\theta'}(\theta)) - k_\theta$ .

The definition of a type game allows us to apply notions and results from standard evolutionary game theory, where evolution acts upon strategies, to the present setting where evolution acts upon types. A similar methodology was used in Mohlin (2012). Note that each type distribution with support in  $C(\mu)$  is represented by a mixed strategy in  $\Gamma_{(\mu, b)}$ .

We want to capture robustness with respect to small groups of individuals, henceforth called *mutants*, which introduce new types and new behaviours into the population. Suppose that a fraction  $\varepsilon$  of the population is replaced by mutants and suppose that the distribution of types within the group of mutants is  $\mu' \in \Delta(\Theta)$ .

Consequently the post-entry type distribution is  $\tilde{\mu} = (1 - \varepsilon) \cdot \mu + \varepsilon \cdot \mu'$ . That is, for each type  $\theta \in C(\mu) \cup C(\mu')$ ,  $\tilde{\mu}(\theta) = (1 - \varepsilon) \cdot \mu(\theta) + \varepsilon \cdot \mu'(\theta)$ . In line with most of the literature on the indirect evolutionary approach we assume that adjustment of behaviour is infinitely faster than the adjustment of the type distribution.<sup>9</sup> Thus we assume that the post-entry type distribution quickly stabilizes into a configuration  $(\tilde{\mu}, \tilde{b})$ . There may exist many such post-entry type configurations, all with the same type distribution, but with different behaviour policies. We note that incumbents do not have to adjust their behaviour against other incumbents in order to continue playing Nash equilibria, and deception equilibria, among themselves. For this reason, we assume (similar to Dekel, Ely, and Yilankaya, 2007) that the incumbents maintain the same pre-entry behaviour among themselves. Formally:

**Definition 5.** Let  $(\mu, b)$  and  $(\tilde{\mu}, \tilde{b})$  be two configurations such that  $C(\mu) \subseteq C(\tilde{\mu})$ . We say that  $(\tilde{\mu}, \tilde{b})$  is *focal* (with respect to  $(\mu, b)$ ) if  $\theta, \theta' \in C(\mu)$  implies that  $\tilde{b}_\theta(\theta') = b_\theta(\theta')$ .

Standard fixed point arguments imply that for every configuration  $(\mu, b)$  and every type distribution  $\tilde{\mu}$  satisfying  $C(\mu) \subseteq C(\tilde{\mu})$ , there exists a behaviour policy  $\tilde{b}$  such that  $(\tilde{\mu}, \tilde{b})$  is a focal configuration.

Our stability notion requires that the incumbents outperform all mutants in all configurations that are focal relative to the initial configuration.

**Definition 6.** A configuration  $(\mu, b)$  is a *neutrally stable configuration (NSC)*, if for every  $\mu' \in \Delta(\Theta)$ , there is some  $\bar{\varepsilon} \in (0, 1)$  such that for all  $\varepsilon \in (0, \bar{\varepsilon})$ , it holds that if  $(\tilde{\mu}, \tilde{b})$ , where  $\tilde{\mu} = (1 - \varepsilon) \cdot \mu + \varepsilon \cdot \mu'$ , is a focal configuration, then  $\mu$  is an NSS in the type game  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ . The configuration  $(\mu, b)$  is an *evolutionarily stable configuration (ESC)* if the same conditions imply that  $\mu$  is an ESS in the type game  $\Gamma_{(\tilde{\mu}, \tilde{b})}$  for each  $\mu' \neq \mu$ .

We conclude this section by discussing four issues related to our notion of stability.

1. The main stability notion that we use in the paper is NSC. The stronger notion of ESC is not useful in our main model because there always exist equivalent types that have slightly different preferences (as the set of preferences is a continuum) and induce the same behaviour as the incumbents. Such mutants would always achieve the same fitness as the incumbents in post-entry configurations, and thus ESCs will never exist. Note that the stability notions in Dekel, Ely, and Yilankaya (2007) and Alger and Weibull (2013) are also based on neutral stability.<sup>10</sup> In Section 6 we study a variant of the model in which the preferences may depend also on the opponent's types. This will allow for the existence of ESCs.
2. Observe that Def. 6 implies internal stability with respect to small perturbations in the frequencies of the incumbent types (because when  $\mu' = \mu$ , then  $\mu$  is required to be an NSS in  $\Gamma_{(\mu, b)}$ ). By standard arguments, internal stability implies that any NSC is “balanced”: all incumbent types obtain the same fitness.
3. By simple adaptations of existing results in the literature, one can show that NSCs and ESCs are dynamically stable. NSCs are Lyapunov stable: no small change in the population composition can lead it away from  $\mu$  in the type game  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ , if types evolve according to the replicator dynamic (Thomas, 1985, Bomze and Weibull, 1995). ESCs are also asymptotically stable: populations starting close enough to  $\mu$  eventually converge to  $\mu$  in  $\Gamma_{(\tilde{\mu}, \tilde{b})}$  if types evolve according to a smooth payoff-monotonic selection dynamic (Taylor and Jonker, 1978, Cressman, 1997, Sandholm, 2010).

<sup>9</sup>Sandholm (2001) and Mohlin (2010) are exceptions.

<sup>10</sup>In their stability analysis of *homo hamiltonensis* preferences Alger and Weibull (2013) disregard mutants who are behaviourally indistinguishable from *homo hamiltonensis* upon entry.

4. The stability notions of Dekel, Ely, and Yilankaya (2007) and Alger and Weibull (2013) only consider monomorphic groups of mutants (i.e., all mutants having the same type). We also consider stability against polymorphic groups of mutants (as do Herold and Kuzmics, 2009). One advantage of our approach is that it allows us to use an adaptation of the well-known notion of ESS, which immediately implies dynamic stability and internal stability, whereas Dekel, Ely, and Yilankaya (2007) have to introduce a novel notion of stability without these properties. We note that our results remain similar with an analogous notion of stability that deals only with monomorphic mutants, except that in this case stability of pure outcomes would imply only a weaker notion of efficiency that compares the fitness only to symmetric profiles, as discussed in Remark 4.2 below.

## 4 Results

### 4.1 Preliminaries

We say that a strategy profile is *efficient* if it maximizes the sum of fitness payoffs. Formally:

**Definition 7.** A strategy profile  $(\sigma, \sigma')$  is *efficient in the game*  $G = (A, \pi)$  if  $\pi(\sigma, \sigma') + \pi(\sigma', \sigma) \geq \pi(a, a') + \pi(a', a)$ , for each action profile  $(a, a')$ .

Let  $\bar{\pi} = \max_{a, a' \in A} (0.5 \cdot (\pi(a, a') + \pi(a', a)))$  denote the efficient payoff, i.e., the average payoff achieved by players who play an efficient profile. A pure Nash equilibrium  $(a, a)$  is strict if  $\pi(a, a) > \pi(a', a)$  for all  $a' \in A$ . Given a configuration  $(\mu, b)$  let  $\bar{n} = \max_{\theta \in C(\mu)} n_\theta$  denote the maximal cognitive level of the incumbents. We refer to incumbents with this cognitive level as the *highest types*.

A deception equilibrium is *fitness maximizing* if it maximizes the fitness of the type with the higher level of the types in the match (under the restriction that the type with the lower level plays an action that is not dominated, given her preferences). Formally:

**Definition 8.** Let  $\theta, \theta'$  be two types with  $n_\theta > n_{\theta'}$ . A deception equilibrium  $(\tilde{\sigma}, \tilde{\sigma}') \in DE(\theta, \theta')$  is *fitness maximizing* if:

$$(\tilde{\sigma}, \tilde{\sigma}') \in \arg \max_{\sigma \in \Delta(A), \sigma' \in \Sigma(u_{\theta'})} \pi(\sigma, \sigma').$$

Let  $FMDE(\theta, \theta') \subseteq DE(\theta, \theta')$  denote the set of all such fitness-maximising deception equilibria of two types  $\theta, \theta'$  with  $n_\theta > n_{\theta'}$ . Note that  $FMDE(\theta, \theta')$  might be an empty set (if there is no action profile that maximizes both the fitness and the subjective utility of the higher type).

A configuration is pure if everyone plays the same action. Formally:

**Definition 9.** A configuration  $(\mu, b)$  is *pure* if there exists  $a^* \in A$  such that  $b_\theta(\theta') = a^*$  for each  $\theta, \theta' \in C(\mu)$ .

With a slight abuse of notation we denote such a pure configuration by  $(\mu, a^*)$ , and we refer to  $a^*$  as the *outcome* of the configuration.

In order to simplify the notation and the arguments in the proofs, we assume throughout this section that the underlying game admits at least three actions (i.e.,  $|A| \geq 3$ ). The results could be extended to games with two actions, but it would make the notation more cumbersome and the proofs less instructive.

### 4.2 Characterization of the Highest Types' Behaviour

In this section we characterize the behaviour of an incumbent type,  $\bar{\theta} = (u, \bar{n})$ , which has the highest level of cognition in the population. We show that the behaviour satisfies the following three conditions:

1. Type  $\bar{\theta}$  plays an efficient action profile when meeting itself.
2. Type  $\bar{\theta}$  maximizes its fitness in all interactions with types with lower cognitive levels.
3. Any opponent with a lower cognitive level achieves at most  $\bar{\pi}$  when being matched with  $\bar{\theta}$ .

**Theorem 1.** *Let  $(\mu^*, b^*)$  be an NSC, and  $\underline{\theta}, \bar{\theta} \in C(\mu^*)$ .*

1. *If  $n_{\bar{\theta}} = \bar{n}$  then  $\pi(b_{\bar{\theta}}(\bar{\theta}), b_{\bar{\theta}}(\bar{\theta})) = \bar{\pi}$ .*
2. *If  $n_{\underline{\theta}} < n_{\bar{\theta}} = \bar{n}$  then  $((b_{\bar{\theta}}(\underline{\theta}), b_{\underline{\theta}}(\bar{\theta}))) \in FMDE(\bar{\theta}, \underline{\theta})$ .*
3. *If  $n_{\underline{\theta}} < n_{\bar{\theta}} = \bar{n}$  then  $\pi(b_{\underline{\theta}}(\bar{\theta}), b_{\bar{\theta}}(\underline{\theta})) \leq \bar{\pi}$ .*

*Proof Sketch (formal proof in Appendix A.2).* The proof utilizes mutants (denoted by  $\theta_1, \theta_2$ , and  $\hat{\theta}$ , below) with the highest cognitive level  $\bar{n}$  and with a specific kind of utility functions, called *indifferent and pro-generous*, which make a player indifferent between all her own actions, but which make the player prefer the opponent to choose an action that allows the player to obtain the highest possible fitness payoff.

To prove part 1 of the theorem, assume to the contrary that  $\pi(b_{\bar{\theta}}(\bar{\theta}), b_{\bar{\theta}}(\bar{\theta})) < \bar{\pi}$ . Let  $a_1, a_2 \in A$  be any two actions such that  $(a_1, a_2)$  is an efficient action profile (i.e.,  $0.5 \cdot (\pi(a_1, a_2) + \pi(a_2, a_1)) = \bar{\pi}$ ). Consider two mutant types  $\theta_1$  and  $\theta_2$ , ( $\theta_1 \neq \theta_2$ ) which are of the highest cognitive level that is present in the population, and have indifferent and pro-generous utility functions. Suppose equal fractions of these two mutant types enter the population. There is a focal post-entry configuration, in which; the incumbents keep playing their pre-entry play among themselves, the mutants play fitness-maximising deception equilibria against lower types, the mutants mimic the play of  $\bar{\theta}$  against all incumbents of level  $\bar{n}$  (and the incumbents behave against the mutants in the same way they behave against  $\bar{\theta}$ ), the mutants mimic  $\bar{\theta}$  when facing a mutant of the same type, and mutants of type  $\theta_1$  play the efficient profile  $(a_1, a_2)$  when they meet mutants of type  $\theta_2$ . In such a focal post-entry configuration mutants on average earn a weakly higher fitness than  $\bar{\theta}$  against the incumbents, and a strictly higher fitness against the mutants. This implies that  $(\mu^*, b^*)$  cannot be an NSC.

To prove part 2, assume to the contrary that  $((b_{\bar{\theta}}(\underline{\theta}), b_{\underline{\theta}}(\bar{\theta}))) \notin FMDE(\bar{\theta}, \underline{\theta})$ . Suppose mutants of type  $\hat{\theta}$  enter. Consider a post-entry configuration in which the incumbents keep playing their pre-entry play among themselves, and the mutants mimic the play of  $\bar{\theta}$ , except that they play a fitness-maximising deception equilibria against all lower types. The mutants obtains a weakly higher payoff than  $\bar{\theta}$  against all types, and a strictly higher payoff than  $\bar{\theta}$  against some lower types. Thus  $(\mu^*, b^*)$  cannot be an NSC.

To prove part 3, assume to the contrary that  $\pi(b_{\underline{\theta}}(\bar{\theta}), b_{\bar{\theta}}(\underline{\theta})) > \bar{\pi}$ . Suppose mutants of type  $\hat{\theta}$  enter. Consider a post-entry configuration in which the incumbents keep playing their pre-entry play among themselves, while the mutants; (i) play fitness-maximising deception equilibria against lower types, (ii) mimic type  $\underline{\theta}$  against type  $\bar{\theta}$ , and (iii) mimic the play of  $\bar{\theta}$  in all other interactions. The type  $\hat{\theta}$  mutants earn strictly more than  $\bar{\theta}$  against both  $\hat{\theta}$  and  $\bar{\theta}$ . The mutants earn weakly more than  $\bar{\theta}$  against all other types. This implies that  $(\mu^*, b^*)$  cannot be an NSC.  $\square$

*Remark.* The first part of Theorem 1 (a highest type must play an efficient strategy when meeting itself) is similar to Dekel, Ely, and Yilankaya's (2007) Proposition 2, which shows that only efficient outcomes can be stable in a setup with perfect observability and no deception. We should note that Dekel, Ely, and Yilankaya (2007) use a weaker notion of efficiency. An action is efficient in the sense of Dekel, Ely, and Yilankaya (2007) (DEY-efficient) if its fitness is highest among the symmetric strategy profiles (i.e., action  $a$  is DEY-efficient if  $\pi(a, a) \geq \pi(\sigma, \sigma)$  for all strategies  $\sigma \in \Delta(A)$ ). Observe that our notion of efficiency (Definition 7) implies

DEY-efficiency, but the converse is not necessarily true. The weaker notion of DEY-efficiency is the relevant one in the set up of Dekel, Ely, and Yilankaya (2007), because they consider only monomorphic groups mutants; i.e., all mutants who enter at the same time are of the same type. A similar result would also hold in our setup, if we imposed a similar limitation on the set of feasible mutants. However, without such a limitation, heterogeneous mutants can correlate their play, and our stronger notion of efficiency is required to characterize stability.

An immediate corollary of Theorem 1 is that a game that has only efficient asymmetric profiles does not admit any NSCs.

**Corollary 1.** *If  $G$  does not have an efficient profile that is symmetric (i.e., if  $\pi(a, a) < \bar{\pi}$  for each  $a \in A$ ), then the game does not admit an NSC.*

### 4.3 Characterization of Pure NSCs

We now present three results which, together with Theorem 1, allow us to completely characterize NSCs with pure outcomes. The first proposition shows that in a pure NSC all incumbents have the minimal cognitive level, since having a higher ability does not yield any advantage when everyone plays the same action.

**Proposition 1.** *If  $(\mu, a^*)$  is an NSC, and  $(u, n) \in C(\mu)$ , then  $n = 1$ .*

*Proof.* Since all players earn the same game payoff of  $\pi(a^*, a^*)$ , they must also incur the same cognitive cost, or else the fitness of the different incumbent types would not be balanced (which contradicts  $(\mu, a^*)$  being an NSC). Moreover, this uniform cognitive level must be level 1. Otherwise a mutant of a lower level, who strictly prefers to play  $a^*$  against all actions, would strictly outperform the incumbents in nearby post-entry focal configurations.  $\square$

The following proposition shows that if  $k_2$  (the cost of having cognitive level 2) is sufficiently small, then any outcome of a pure NSC must be a Nash equilibrium of the underlying game. The reason is that if the pure outcome were not a Nash equilibrium, then the population could be invaded by mutants of cognitive level 2, who deceive the incumbents into thinking they face other incumbents, and best-reply to the incumbents' play.

**Proposition 2.** *Suppose*

$$k_2 < \bar{k} := \min_{a, a', a'' \text{ s.t. } \pi(a, a'') \neq \pi(a', a'')} |\pi(a, a'') - \pi(a', a'')|. \quad (1)$$

*If  $(\mu, a^*)$  is an NSC, then  $(a^*, a^*)$  is a symmetric Nash equilibrium, in fitness payoffs.*

*Proof.* Assume to the contrary that  $(\mu, a^*)$  is a pure NSC and  $a^*$  is not a best response to itself; i.e., there exist  $a' \in A$  such that  $\pi(a', a^*) > \pi(a^*, a^*)$ . Assume without loss of generality that  $a'$  is a best reply against  $a^*$  (in fitness terms). By Proposition 1, all incumbents have cognitive level 1. Consider a mutant  $\theta' = (\pi, 2)$  of cognitive level 2 and materialistic preferences. There is a focal post-entry configuration in which mutants play the deception equilibrium  $(a', a^*)$  against the incumbents. Observe that the payoff is strictly higher when mutants face an incumbent than when two incumbents face each other:

$$\pi(a', a^*) - k_2 > \pi(a', a^*) - \bar{k} \geq \pi(a^*, a^*).$$

This implies that if the mutants are sufficiently rare, they outperform the incumbents in the post-entry focal configuration.  $\square$

The next proposition shows that any action that is both efficient and a strict Nash equilibrium, can be induced as the outcome of an NSC. The intuition is as follows (and is similar to Dekel, Ely, and Yilankaya, 2007, Proposition 6). Consider a monomorphic population in which all individuals are of cognitive level 1 and have a dominant action that is an efficient strict Nash equilibrium action. The fact that the action profile is strict Nash and efficient, implies that any group of mutants is weakly outperformed.

**Proposition 3.** *If  $(\bar{a}, \bar{a})$  is both efficient and a strict Nash equilibrium (in fitness payoffs), then there exists a type distribution  $\mu$  such that  $(\mu, \bar{a})$  is an NSC.*

*Proof.* Consider a monomorphic configuration  $(\mu, \bar{a})$  consisting of type  $(\theta^*, 1)$  where all incumbents are of cognitive level 1 and of the same preference type  $\theta^*$ , which strictly prefers to play  $\bar{a}$  regardless of what the opponent plays. Observe, that after any mutant's entry, in all focal post-entry configurations the incumbent  $\theta^*$  will always play  $\bar{a}$  (since  $\bar{a}$  is strictly dominant for  $\theta^*$ ). Since the incumbent is always playing  $\bar{a}$ , and  $(\bar{a}, \bar{a})$  is a strict Nash equilibrium of  $G$ , mutants who do not play  $\bar{a}$  when they are matched with  $\theta^*$  will obtain strictly less fitness than the incumbents if their population share is sufficiently small. But for mutants who play  $\bar{a}$  whenever they are matched with  $\theta^*$ , the incumbents' average fitness is given by  $\pi(\bar{a}, \bar{a})$ , and since mutants cannot obtain an average fitness strictly higher than this when they are matched among themselves (since  $(\bar{a}, \bar{a})$  is efficient), they cannot obtain a strictly higher average fitness either. We conclude that  $(\mu, \bar{a})$  is an NSC.  $\square$

Part 1 of Theorem 1, and with Propositions 1–3 imply that pure NSCs outcomes are essentially efficient Nash equilibria. Formally:

**Corollary 2 (Characterization of pure NSCs).**

1. *If  $(a^*, a^*)$  is both efficient and a strict Nash equilibrium in fitness payoffs, then it is the outcome of a pure NSC.*
2. *If  $(a^*, a^*)$  is the outcome of a pure NSC and  $k_2 < \bar{k}$ , then it is both efficient and a Nash equilibrium in fitness payoffs.*

#### 4.4 Characterization of NSCs in Generic Games

In this section we characterize NSCs in generic games, by which we mean games in which any two different action profiles both give an individual player different payoffs, and yield different total payoffs.

**Definition 10.** A (symmetric) game is generic if for each  $a, a', b, b' \in A$ ,  $\{a, a'\} \neq \{b, b'\}$  implies:

$$\pi(a, a') \neq \pi(b, b'), \text{ and } \pi(a, a') + \pi(a', a) \neq \pi(b, b') + \pi(b', b).$$

For example, if the entries of the payoff matrix  $\pi$  are drawn from a continuous distribution on an open subset of the real numbers, then the induced game is generic with probability one.

Note that a generic game admits at most one efficient action profile. From Corollary 1 we know that if the game does not have a symmetric efficient profile then it does not admit any NSC. Hence we can restrict attention to games with exactly one efficient action. Let  $\bar{a}$  denote this unique efficient action.

Next we present the main result of the paper: all incumbent types play efficiently in any NSC of a generic game.

**Theorem 2.** *If  $(\mu^*, b^*)$  is an NSC of a generic game with a (unique) efficient outcome  $(\bar{a}, \bar{a})$ , then  $(b_\theta(\theta'), b_{\theta'}(\theta)) = (\bar{a}, \bar{a})$ , for all  $\theta, \theta' \in C(\mu^*)$ .*

*Proof.* Assume to the contrary that configuration  $(\mu^*, b^*)$  is an NSC such that  $(b_\theta(\theta'), b_{\theta'}(\theta)) \neq (\bar{a}, \bar{a})$  for some types  $\theta, \theta' \in C(\mu^*)$ . Let  $\hat{\theta}$  be the type with the highest cognitive level among the types that satisfy at least one of the following conditions:

- (A)  $\hat{\theta}$  plays inefficiently against itself, i.e.,  $(b_{\hat{\theta}}(\hat{\theta}), b_{\hat{\theta}}(\hat{\theta})) \neq (\bar{a}, \bar{a})$ .
- (B)  $\hat{\theta}$  plays inefficiently against a weakly higher type, i.e.,  $(b_{\hat{\theta}}(\theta'), b_{\theta'}(\hat{\theta})) \neq (\bar{a}, \bar{a})$  for some  $\theta' \neq \hat{\theta}$  with  $n_{\hat{\theta}} \leq n_{\theta'}$ .
- (C) A strictly lower type earns strictly more than  $\bar{\pi}$  against  $\hat{\theta}$ , i.e.,  $\pi(b_{\hat{\theta}}(\theta''), b_{\theta''}(\hat{\theta})) > \bar{\pi}$  for some  $\theta'' \neq \hat{\theta}$  with  $n_{\hat{\theta}} > n_{\theta''}$ .

We will now successively rule out each of these cases.

Assume first that (A) holds. Let  $\hat{u}$  be a utility function that is identical to  $u_\theta$  except that: (i) the payoff of the outcome  $(\bar{a}, \bar{a})$  is increased by the minimal amount required to make it a best reply to itself, and (ii) the payoff of some other outcome is altered slightly (to ensure  $\hat{u}$  is not already an incumbent) in a way that does not force  $\hat{\theta}$  to behave differently from  $\theta$ . (The formal definition of  $\hat{u}$  is provided in Appendix A.3.) Suppose that a mutant of type  $\hat{\theta} = (\hat{\theta}, n_\theta)$  enters. Consider a focal post-entry configuration in which the type  $\hat{\theta}$  mutants mimic the play of the type  $\hat{\theta}$  incumbents in all matches except: (i) the mutants play the efficient profile  $(\bar{a}, \bar{a})$  among themselves (which yields a higher payoff than what  $\hat{\theta}$  achieves when matched against  $\hat{\theta}$ ), and (ii) when the mutants face a higher type they play either  $(\bar{a}, \bar{a})$  or the same deception equilibrium that the higher types play against  $\hat{\theta}$ . It follows that the mutants  $\hat{\theta}$  earn a strictly higher payoff than  $\hat{\theta}$  against  $\hat{\theta}$ , and a weakly higher fitness than type  $\theta$  against all other types. Thus the mutants strictly outperform the incumbents, which contradicts the assumption that  $(\mu^*, b^*)$  is an NSC. The full technical details of this argument are given in Appendix A.3.

Next, assume that case (B) holds and that case (A) does not hold. This implies that:

$$0.5 \cdot (\pi(b_\theta(\theta'), b_{\theta'}(\theta)) + \pi(b_{\theta'}(\theta), b_\theta(\theta'))) < \bar{\pi} = \pi(b_\theta(\theta), b_\theta(\theta)) = \pi(b_{\theta'}(\theta'), b_{\theta'}(\theta')).$$

That is, in the subpopulation that includes types  $\theta$  and  $\theta'$  the within-type matchings yield higher payoffs than out-group matchings (an average payoff of less than  $\bar{\pi}$ ). The following formal argument shows that this property implies dynamic instability. The fact that  $(\mu^*, b^*)$  is an NSC implies that  $\mu^*$  is an NSS in the type game  $\Gamma_{(\mu^*, b^*)}$ . Let  $B$  be the payoff matrix of the type game  $\Gamma_{(\mu^*, b^*)}$  and let  $n = |C(\mu^*)|$ . It is well known (e.g., Hofbauer and Sigmund, 1988, Exercise 6.4.3, and Hofbauer, 2011, pp. 1–2) that an interior Nash equilibrium of a normal form game can be an NSS if and only if the payoff matrix is negative semi-definite with respect to the tangent space, i.e., if and only if  $x \cdot Bx \leq 0$  for each  $x \in \mathbb{R}^n$  such that  $\sum_i x_i = 0$ . Assume without loss of generality that type  $\theta$  ( $\theta'$ ) is represented by the  $j^{th}$  ( $k^{th}$ ) row of the matrix  $B$ . Let the column vector  $x$  be defined as follows:

$$x(i) = \begin{cases} 1 & i = j \\ -1 & i = k \\ 0 & i \neq j, k. \end{cases}$$

That is, the vector  $x$  has all entries equal to zero, except the  $j^{th}$  entry which is equal to 1, and the  $k^{th}$  entry

which is equal to  $-1$ . We have:

$$\begin{aligned}
x \cdot Bx &= B_{jj} - B_{jk} - B_{jk} + B_{kk} \\
&= \bar{a} - k_{n_{\theta_j}} + \bar{a} - k_{n_{\theta_k}} - \left( \pi(b_{\theta'}(\theta'), b_{\theta'}(\theta)) - k_{n_{\theta}} + \pi(b_{\theta'}(\theta), b_{\theta'}(\theta')) - k_{n_{\theta_k}} \right) \\
&= 2 \cdot \bar{a} - \left( \pi(b_{\theta_j}(\theta_k), b_{\theta_k}(\theta_j)) + \pi(b_{\theta_k}(\theta_j), b_{\theta_j}(\theta_k)) \right) > 0.
\end{aligned}$$

Thus  $B$  is not negative semidefinite.

Finally, assume that only case (C) holds. Let  $\bar{\theta}$  be an incumbent type with the highest cognitive level. The fact that case (B) does not hold implies that  $(b_{\bar{\theta}}(\bar{\theta}), b_{\bar{\theta}}(\hat{\theta})) = (\bar{a}, \bar{a})$ , and so  $\pi(b_{\bar{\theta}}(\hat{\theta}), b_{\bar{\theta}}(\bar{\theta})) = \bar{\pi}$ . The fact that case (C) holds implies that  $\pi(b_{\theta''}(\hat{\theta}), b_{\bar{\theta}}(\theta'')) > \bar{\pi}$  (and, in particular, that  $b_{\bar{\theta}}(\theta'') \in \Sigma_{\hat{\theta}}$ ). This contradicts part (2) of Theorem 1, according to which we should have  $(b_{\bar{\theta}}(\hat{\theta}), b_{\bar{\theta}}(\bar{\theta})) = FMDE(\bar{\theta}, \hat{\theta})$ .  $\square$

Note that in a generic game any pure action is either a strict equilibrium, or not a best reply to itself. Combining the results of this section with the above characterization of pure NSCs yields the following corollary, which fully characterizes the NSCs of generic games.

**Corollary 3 (Characterization of NSCs in Generic Games).** *Let  $G$  be a generic game.*

1. *If  $(a^*, a^*)$  is an efficient Nash equilibrium in fitness payoffs, then  $(a^*, a^*)$  is the outcome of a pure NSC.*
2. *If  $(\mu^*, b^*)$  is an NSC and  $k_2 < \bar{k}$ , then the NSC is pure, i.e.,  $b^* \equiv a^*$ , with the outcome  $(a^*, a^*)$  being an efficient Nash equilibrium in fitness payoffs.*
3. *If  $G$  does not have an efficient symmetric Nash equilibrium, then  $G$  does not admit any NSC.*

## 4.5 Non-Pure NSCs in Non-generic Games

The previous two subsections fully characterizes (i) pure NSCs and (ii) NSCs in generic games. In this section we analyse non-pure NCSs in non-generic games. Non-generic games may be of interest in various setups, such as: (1) normal-form representation of generic extensive form games (the induced matrix is typically non-generic), and (2) special interesting families of games, such as zero-sum games. Unlike generic games, non-generic games can admit (non-pure) NSCs with multiple cognitive levels and non-Nash behaviour. To demonstrate this we consider the Rock-Paper-Scissors game, with the following payoff matrix:

	$R$	$P$	$S$
$R$	0, 0	-1, 1	1, -1
$P$	1, -1	0, 0	-1, 1
$S$	-1, 1	1, -1	0, 0

The result below shows that, under mild assumptions on the cognitive costs function, this game admits an NSC in which all players have the same materialistic preferences, but players of different cognitive levels coexist, and non-Nash profiles are played in all matches between two individuals of different cognitive levels. More precisely, when individuals of different cognitive levels meet, the higher-level individual deceives the lower-level individual to take a pure action that the higher-level individual then best-responds to. Thus the higher-level individual earns 1 and her opponent earns  $-1$ . Individuals of the same cognitive level play the unique Nash equilibrium. This means that higher-level types will obtain a payoff of 1 more often than lower-level types,



and lower-level types will obtain a payoff of  $-1$  more often than higher-level types. In the NSC this payoff difference is offset exactly by the higher cognitive cost paid by the higher types. Moreover, the cognitive cost is increasing such that at some point the cost of cognition outweighs any payoff differences that may arise from the underlying game. This implies that there is an upper bound on the cognitive sophistication in the population.

**Proposition 4.** *Let  $G$  be a Rock-Paper-Scissors game. Let  $u^\pi$  denote the (materialistic) preference such that  $u^\pi(a, a') = \pi(a, a')$  for all profiles  $(a, a')$ . Suppose that the marginal cognitive cost is small but non-vanishing, so that (a) there is an  $N$  such that  $k_N \leq 2 < k_{N+1}$ , and (b) it holds that<sup>11</sup>  $1 > k_{n+1} - k_n$  for all  $n \leq N$ . There exists an NSC  $(\mu^*, b^*)$ , such that  $C(\mu^*) \subseteq \{(u^\pi, n)\}_{n=1}^N$ , and  $\mu^*$  is mixed (i.e.,  $|C(\mu^*)| > 1$ ). The behaviour of the incumbent types is as follows: if the individuals in a match are of different cognitive levels, then the higher level plays Paper and the lower level plays Rock; if both individuals in a match are of the same cognitive level, then they both play the unique Nash equilibrium (i.e., randomize uniformly over the three actions).*

Appendix C.2 contains a formal proof of this result and relates it to a similar construction in Conlisk (2001).

## 5 Non-observability among Cognitive Equals

The main model assumes that types with equal cognitive levels can observe each others preferences. In this section we present an alternative “tie-breaking rule” according to which a player is unable to observe the preferences of an opponent with the same cognitive level.

### 5.1 Changes to the Baseline Model

**Play between Types with Equal Cognitive Levels.** If two individuals of the same cognitive level are matched to play, then they play a Bayes–Nash equilibrium of the Bayesian game in which each player knows only that her opponent has the same cognitive level (but she cannot observe her opponent’s preferences). Given a distribution of types  $\mu$ , we say that  $n \in \text{proj}_N C(\mu)$  if there exists a type  $(u, n) \in C(\mu)$ . Given a distribution  $\mu$  and  $n \in \text{proj}_N C(\mu)$ , let  $\Theta_n$  denote the set of types with level  $n$ , and let  $\mu_n$  denote the distribution of types conditional on having cognitive level  $n$ :

$$\mu_n(u, n) = \frac{\mu(u, n)}{\sum_{(u', n) \in C(\mu)} \mu(u', n)}.$$

Given distribution  $\mu$  and  $n \in \text{proj}_N C(\mu)$ , a vector of action profiles  $(\sigma_\theta)_{\theta \in \Theta_n}$  is a Bayes–Nash equilibrium of the game between types with level  $n$ , if each player best-responds to the aggregate behaviour of players of the same cognitive level. Let  $BNE(n) \subseteq (\Delta A)^{|\Theta_n|}$  be the set of Bayes–Nash equilibria of the game induced by the interactions between players of cognitive level  $n$ , i.e.,

$$BNE(n) = \left\{ (\sigma_\theta)_{\theta \in \Theta_n} \in (\Delta A)^{|\Theta_n|} : \forall \sigma_{(u, n)} \in \Theta_n, \sigma_{(u, n)} \in \underset{\sigma \in \Delta(A)}{\text{argmax}} \sum_{\theta \in \Theta_n} \mu_n(\theta) \cdot u(\sigma, \sigma_\theta) \right\}.$$

The play among types with equal levels is analogous to the play of all types in the unobservable case of Dekel, Ely, and Yilankaya (2007, Section 4). Players with different cognitive levels play deception equilibria as in the baseline model.

<sup>11</sup>If we define  $\bar{k}$  as in Eq. (1) then we have  $\delta = 1$ , in Proposition 4. Thus the condition that  $\bar{k} > k_{n+1} - k_n$  for all  $n$  in Proposition 4, may be viewed as an extension of the condition  $k_2 < \bar{k}$  in Proposition 2.

**Redefining Configurations.** We redefine the notion of a configuration as follows:

**Definition 11.** A *configuration* is a pair  $(\mu, b)$  where  $\mu \in \Delta(U)$  is a type distribution, and  $b : C(\mu) \times C(\mu) \rightarrow \Delta(A)$  is a *behaviour policy* such that:

1. Players play the same against all opponents of the same level, i.e., for each  $n \in \text{proj}_N C(\mu)$ , if  $\theta, \theta', \theta'' \in \Theta_n$  then  $b_\theta(\theta') = b_\theta(\theta'')$ .
2. Types with the same level play a Bayes–Nash equilibrium, i.e., for each  $n \in \text{proj}_N C(\mu)$

$$((b_\theta(\theta))_{\theta \in \Theta_n}) \in BNE(n).$$

3. Types with different levels play deception equilibria, i.e., for each  $\theta, \theta' \in C(\mu)$  :

$$n_\theta > n_{\theta'} \implies (b_\theta(\theta'), b_{\theta'}(\theta)) \in DE(\theta, \theta').$$

**Redefining NSC.** Unlike in the baseline model, focal configurations (in which the incumbents keep their pre-entry play) do not always exist. Accordingly, we modify the definition of NSC slightly to require the existence of a focal configuration for any given distribution of mutants.<sup>12</sup> Formally:

**Definition 12.** A configuration  $(\mu, b)$  is a *neutrally stable configuration (NSC)*, if for every  $\mu' \in \Delta(\Theta)$ , there is some  $\bar{\varepsilon} \in (0, 1)$  such that for all  $\varepsilon \in (0, \bar{\varepsilon})$ , it holds that; (1) there exists a focal configuration  $(\tilde{\mu}, \tilde{b})$ , where  $\tilde{\mu} = (1 - \varepsilon) \cdot \mu + \varepsilon \cdot \mu'$ , and (2) for each such focal configuration  $(\tilde{\mu}, \tilde{b})$ , the distribution  $\mu$  is an NSS in the type game  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ .

## 5.2 Results

The following simple results show that with non-observable equals, it is no longer the case that stable outcomes must be efficient. The first result shows that any strict Nash equilibrium (even an inefficient) is the outcome of a pure NSC. The argument is similar to Dekel, Ely, and Yilankaya (2007, Proposition 5-b) and is presented briefly.

**Proposition 5.** *If  $(a^*, a^*)$  is a strict Nash equilibrium, then there exists a pure NSC  $(\mu, a^*)$ , where  $\mu$  contains a single type  $(1, u^*)$  and  $a^*$  is a strictly dominant action given the utility function  $u^*$ .*

*Proof.* Observe that the incumbents always play  $a^*$  in any focal configuration (and that such a focal configuration always exists). The fact that  $a^*$  is a strict equilibrium implies that any mutant who does not always play  $a^*$  against the incumbents is strictly outperformed if the mutants are sufficiently rare. If the mutant type has level one, it implies that she must always play  $a^*$  also against mutants, and that she achieves the same payoff as the incumbents. If the mutant has higher cognitive level, then if the mutants are sufficiently rare the higher cognitive cost implies that the mutant is strictly outperformed.  $\square$

The next result shows that players with the highest type must play a Nash equilibrium (of the underlying game) among themselves. The argument is similar to Dekel, Ely, and Yilankaya (2007, Proposition 5-a) and is presented briefly.

---

<sup>12</sup>All the results remain the same if one only requires  $\delta$ -focality à la Dekel, Ely, and Yilankaya (2007)

**Proposition 6.** *Let  $(\mu^*, b^*)$  be an NSC and let  $\bar{n}$  be the highest cognitive level. The action profile  $(b_\theta(\theta))_{\theta \in \Theta_{\bar{n}}}$  is a Bayes–Nash equilibrium of the fitness game. That is:*

$$\forall \theta^* \in \Theta_{\bar{n}}, b_{\theta^*}(\theta^*) \in \operatorname{argmax}_{\sigma \in \Delta(A)} \sum_{\theta \in \Theta_n} \mu_n(\theta) \cdot \pi(\sigma, b_\theta(\theta)).$$

*Proof.* Assume to the contrary that the highest types do not play a Bayes–Nash equilibrium of the fitness game. Let  $\bar{\theta} \in \Theta_{\bar{n}}$  be one such type. That is:

$$b_{\bar{\theta}}(\bar{\theta}) \notin \operatorname{argmax}_{\sigma \in \Delta(A)} \sum_{\theta \in \Theta_n} \mu_n(\theta) \cdot \pi(\sigma, b_\theta(\theta)).$$

Let  $\theta' = (u', \bar{n})$  be a mutant type which is indifferent between any two outcomes, and which mimics the play of  $\theta^*$  against lower types, while best-replying in the induced fitness game between the highest types. Such a mutant type would strictly outperform the incumbent  $\bar{\theta}$  in any focal configuration.  $\square$

## 6 Type-Interdependent Preferences

In this section we describe an extension of our baseline model, such that the preferences may depend not only on action profiles, but also on the opponent’s type.

### 6.1 Changes to the Baseline Model

We briefly describe how to amend the model to handle type-interdependent preferences. Our construction is similar to that of Herold and Kuzmics (2009).

When the preferences of a type depend on the opponent’s type, we can no longer work with the set of all possible preferences, because it would create problems of circularity and cardinality.<sup>13</sup> Instead, we must restrict attention to a pre-specified set of feasible preferences. We begin by defining  $\Theta_{ID}$  as an arbitrary set of labels. Each label is a pair  $\theta = (u, n) \in \Theta_{ID}$ , where  $n \in \mathbb{N}$  and  $u$  is a type-interdependent utility function that depends on the played action profile as well as the opponent’s label,

$$u : A \times A \times \Theta_{ID} \rightarrow \mathbb{R}.$$

Each label  $\theta = (u, n)$  may now be interpreted as a type. The definition of  $u$  extends to mixed actions in the obvious way. We use the label  $u$  also to describe its associated utility function  $u$ . Thus  $u(\sigma, \sigma', \theta')$  denotes the subjective payoff that a player with preferences  $u$  earns when she plays strategy  $\sigma$  against an opponent with type  $\theta'$  who plays strategy  $\sigma'$ .

Let  $U_{ID}$  denote the set of all preferences that are part of some type in  $\Theta_{ID}$ , i.e.,  $U_{ID} = \{u : \exists n \in \mathbb{N} \text{ s.t. } (u, n) \in \Theta_{ID}\}$ . For each type-neutral preference  $u \in U$  we can define an equivalent type-interdependent preference  $u \in U_{ID}$ , which is independent of the opponent’s type; that is,  $u'(\sigma, \sigma', \theta') = u''(\sigma, \sigma', \theta'')$  for each  $u', u'' \in U_{ID}$ . Let  $U_N$  denote the set of all such type-interdependent versions of the type-neutral preferences of the baseline model. All of our results allow, but do not require, that  $U_N \subseteq U_{ID}$ .

<sup>13</sup>The circularity comes from the fact that each type contains a preferences component, which is identified with a utility function defined over types (and action profiles). To see that this creates a problem if the set of types is unrestricted, let  $\Theta_*$  be the set of types and suppose that the corresponding set of preferences,  $U_*$ , contains all mappings  $u : A \times A \times \Theta_* \rightarrow \mathbb{R}$ . The cardinality of this set is  $|U| \cdot |\Theta_*|$ , but if  $U_*$  is indeed the set of *all* mappings  $u : A \times A \times \Theta_* \rightarrow \mathbb{R}$ , then we must have  $|U_*| = |U| \cdot |\Theta_*|$ . Since  $|\Theta_*| \geq |U_*|$  this is a contradiction. See also footnote 10 in Herold and Kuzmics (2009).

Next, we amend the definitions of Nash equilibrium, undominated strategies, and deception equilibrium. The best-reply correspondence now takes both strategies and types as arguments:  $BR_u(\sigma', \theta') = \arg \max_{\sigma \in \Delta(A)} u(\sigma, \sigma', \theta')$ . Accordingly we adjust the definition of the set of Nash equilibria,

$$NE(\theta, \theta') = \{(\sigma, \sigma') \in \Delta(A) \times \Delta(A) : \sigma \in BR_u(\sigma', \theta') \text{ and } \sigma' \in BR_{u'}(\sigma, \theta)\},$$

and the set of *undominated strategies*

$$\Sigma(\theta) = \{\sigma \in \Delta(A) : \text{there exists } \sigma' \in \Delta(A) \text{ and } \theta' \in \Theta_{ID} \text{ such that } \sigma \in BR_u(\sigma', \theta')\}.$$

Finally, we adapt the definition of deception equilibrium. Given two types  $\theta, \theta'$  with  $n_\theta > n_{\theta'}$ , a strategy profile  $(\tilde{\sigma}, \tilde{\sigma}')$  is a *deception equilibrium* if

$$(\tilde{\sigma}, \tilde{\sigma}') \in \arg \max_{\sigma \in \Delta(A), \sigma' \in \Sigma(\theta')} u_\theta(\sigma, \sigma', \theta').$$

Let  $DE(\theta, \theta')$  be the set of all such deception equilibria. The rest of our model remains unchanged.

## 6.2 Pure Maxmin and Minimal Fitness

The pure maxmin and minmax values give a minimal bound to the fitness of an NSC. Given a game  $G = (A, \pi)$ , define  $\underline{M}(\bar{M})$  as its pure maxmin (minmax) value:

$$\underline{M} = \max_{a_1 \in A} \min_{a_2 \in A} \pi(a_1, a_2), \quad \bar{M} = \min_{a_2 \in A} \max_{a_1 \in A} \pi(a_1, a_2).$$

The pure maxmin value  $\underline{M}$  is the minimal fitness payoff a player can guarantee herself in the sequential game in which she plays first, and the opponent replies in an arbitrary way (i.e., not necessarily in a way that maximizes the opponent's fitness.) The pure minmax value  $\bar{M}$  is the minimal fitness payoff a player can guarantee herself in the sequential game in which her opponent plays first an arbitrary action, and she best-plies to the opponent's pure action. It is immediate that  $\underline{M} \leq \bar{M}$ , and that the minmax value in mixed actions is between these two values.

Let  $a_M$  be a maxmin action of a player; an action  $a_M$  guarantees that the player's payoff is at least  $\underline{M}$ ,

$$a_M \in \arg \max_{a_1 \in A} \min_{a_2 \in A} \pi(a_1, a_2).$$

The next proposition (which holds also in the baseline model with type-neutral preferences) shows that the maxmin value is a lower bound on the fitness payoff obtained in an NSC. The intuition is that if the payoff is lower, then a mutant of cognitive level 1, with preferences such that the maxmin action  $a_M$  is dominant, will outperform the incumbents.

**Definition 13.** Given a pure action  $a^* \in A$ , let  $u^{a^*} \in U_N$  be the (type-neutral) preferences in which the player obtains a payoff of 1 if she plays  $a^*$  and a payoff of 0 otherwise (i.e.,  $a^*$  is a dominant action regardless of the opponent's preferences).

**Proposition 7.** *Suppose that  $(u^{a_M}, 1) \in \Theta_{ID}$ . If  $(\mu, b)$  is an NSC then  $\Pi(\mu, b) \geq \underline{M}$ .*

*Proof.* Assume to the contrary that  $\Pi(\mu, b) < \underline{M}$ . Consider a monomorphic group of mutants with type

$(u^{a_M}, 1)$ . The fact that  $a_M$  is a maxmin action implies that

$$\pi_{(u^{a_M}, 1)}((\tilde{\mu}, \tilde{b})) \geq \underline{M}$$

in any post-entry configuration. Furthermore, due to continuity it holds that  $\Pi_\theta(\tilde{\mu}, \tilde{b}) < \underline{M}$  for any  $\theta \in C(\mu)$  in all sufficiently close focal post-entry configuration. This contradicts  $\mu$  being an NSS in  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ , and thus it contradicts  $(\mu, b)$  being an NSC.  $\square$

### 6.3 characterization of Pure Stable Configurations

In this subsection we show that, essentially, a pure action can be an outcome of an ESC if and only if it is a Nash equilibrium that yields each player a payoff above her minmax/maxmin value.

We first adapt Propositions 1–2 to the current setup. Specifically, we show that if  $(\mu^*, a^*)$  is a pure NSC, then: (1) all incumbents have the same cognitive level, and (2)  $a^*$  is a symmetric Nash equilibrium, provided that the marginal cognitive costs are sufficiently small (smaller than  $\bar{k}$ , as defined in Eq. (1)).

**Proposition 8.** *If  $(\mu^*, a^*)$  is a pure NSC then the following holds:*

1. *If  $\theta, \theta' \in C(\mu^*)$  then  $n_\theta = n_{\theta'}$ .*
2. *If  $k_{n+1} - k_n < \bar{k}$  for each  $n$ , then  $(a^*, a^*)$  is a symmetric Nash equilibrium in fitness payoffs.*

*Proof.*

1. Since all players earn the same game payoff of  $\pi(a^*, a^*)$ , they must also incur the same cognitive cost, or else the fitness of the different incumbent types would not be balanced (which contradicts the fact that  $(\mu, a^*)$  is an NSC).
2. Assume to the contrary that there exist  $a' \in A$  such that  $\pi(a', a^*) > \pi(a^*, a^*)$ . Assume without loss of generality that  $a'$  is a best reply against  $a^*$  (in fitness terms). By the previous part all the incumbents have the same cognitive level  $n$ . Consider a monomorphic group of mutants  $\theta' = (\pi, n+1)$ . There is a focal post-entry configuration in which the mutants play the deception equilibrium  $(a', a^*)$  against the incumbents. Observe that the mutants obtain a strictly higher payoff when facing an incumbent than the payoff of two incumbents who face each other:

$$\pi(a', a^*) - (k_{n+1} - k_n) > \pi(a', a^*) - \bar{k} \geq \pi(a^*, a^*).$$

This implies that if the mutants are sufficiently rare, they outperform the incumbents in the post-entry focal configuration.  $\square$

Let  $a_{\bar{M}}$  be a minmax action, i.e., an action that guarantees that the opponent's payoff is at most  $\bar{M}$ ;

$$a_{\bar{M}} \in \arg \min_{a_2 \in A} \max_{a_1 \in A} \pi(a_1, a_2).$$

**Definition 14.** Given any two actions  $\tilde{a}, \tilde{a}' \in A$ , let  $u_{\tilde{a}}$  be the discriminating preferences defined by the following utility function: For all  $a'$ ,

$$u_{\tilde{a}}(a, a', \theta') = \begin{cases} 1 & \text{if } u_{\theta'} = u_{\tilde{a}} \text{ and } a = \tilde{a} \\ 1 & \text{if } u_{\theta'} \neq u_{\tilde{a}} \text{ and } a = \tilde{a}' \\ 0 & \text{otherwise.} \end{cases}$$

In words, the preferences  $u_{\tilde{a}}$  are such that  $\tilde{a}$  is a dominant action against an opponent with the same preferences, and  $\tilde{a}'$  is the dominant action against all other opponents.

The following result shows that any action  $a^*$  that is both a symmetric Nash equilibrium and yields a payoff above the minmax value can be implemented as the unique pure outcome of an ESC. (Recall that  $\theta$  is used to denote that probability distribution  $\mu$  puts all weight on  $\theta$ , i.e.,  $\mu(\theta) = 1$ .)

**Proposition 9.** *Suppose that  $(u_{a_M}^{a^*}, 1) \in \Theta_{ID}$ . If action  $(a^*, a^*)$  is a symmetric Nash equilibrium and  $\pi(a^*, a^*) > \bar{M}$ , then  $((u_{a_M}^{a^*}, 1), a^*)$  is an ESC.*

*Proof.* Suppose that all incumbents are of type  $(u_{a_M}^{a^*}, 1)$ . Note that in all focal post-entry configurations the incumbent  $(u_{a_M}^{a^*}, 1)$  always plays either  $a^*$  or  $a_M$ . Against a mutant  $(\theta, 1)$  with cognitive level 1, an incumbent plays  $a^*$  if and only if  $u(\theta) = u_{a_M}^{a^*}$ . The fact that  $\pi(a^*, a^*) > \bar{M}$  implies that any mutant  $\theta \neq (u_{a_M}^{a^*}, 1)$  earns a strictly lower payoff against the incumbents in any post-entry configuration. As a result, if the frequency of mutants is sufficiently small, then they are strictly outperformed. Against a mutant  $(\theta, n)$  with cognitive level  $n > 1$ , an incumbent may play either  $a^*$  or  $a_M$ . Since  $a^*$  is a symmetric Nash equilibrium and  $\pi(a^*, a^*) > \bar{M}$  the mutants earn at most  $\pi(a^*, a^*)$  in matches against incumbents. Consequently, as the fraction of mutants vanishes the average fitness of mutants is weakly less than  $\pi(a^*, a^*) - k_n$ , and the average fitness of the incumbents is  $\pi(a^*, a^*)$ . Since  $k$  is strictly increasing this implies that  $((u_{a_M}^{a^*}, 1), a^*)$  is an ESC.  $\square$

The results of this section imply the following corollary, which characterizes pure outcomes of stable configurations in terms of being Nash equilibria that yield payoffs above the pure maximin/minmax values.

**Corollary 4.** *Suppose that  $(u^{a_M}, 1) \in \Theta_{ID}$  and  $(u_{a_M}^{a^*}, 1) \in \Theta_{ID}$ .*

1. *If  $(a^*, a^*)$  is a Nash equilibrium and  $\pi(a^*, a^*) > \bar{M}$ , then  $a^*$  is the outcome of a pure ESC.*
2. *If  $a^*$  is the outcome of a pure NSC, and  $k_{n+1} - k_n < \bar{k}$  for all  $n$ , then  $(a^*, a^*)$  is a symmetric Nash equilibrium and  $\pi(a^*, a^*) \geq \underline{M}$ .*

## 6.4 Application: In-group Cooperation and Out-group Exploitation

The following table represents a family of Hawk-Dove games. When both players play  $D$  (Dove) they earn 1 each and when they both play  $H$  (Hawk) they earn 0. When a player plays  $H$  against an opponent playing  $D$ , she obtains an additional gain of  $g > 0$  and the opponent incurs a loss of  $l \in (0, 1)$ .

	$H$	$D$	
$H$	0, 0	$1 + g, 1 - l$	(2)
$D$	$1 - l, 1 + g$	1, 1	

It is natural to think of mutual play of  $D$  as the cooperative outcome. We define preferences that induce players to cooperate with their own kind and to seek to exploit those who are not of their own kind.

**Definition 15.** Let  $u^n$  denote the preferences such that:

- (1) If  $u_{\theta'} = u^n$  and  $n_{\theta'} = n$  then  $u^n(D, a', \theta') = 1$  and  $u^n(H, a', \theta') = 0$  for all  $a'$ .
- (2) If  $u_{\theta'} \neq u^n$  or  $n_{\theta'} \neq n$  then  $u^n(H, a', \theta') = 1$  and  $u^n(D, a', \theta') = 0$  for all  $a'$ .

Thus, facing someone who is of the same type, an individual with  $u^n$ -preferences strictly prefers cooperation, in the sense of playing  $D$ . When facing someone who is not of the same type, an individual with  $u^n$ -preferences prefers the exploitative outcome  $(H, D)$ , and after that she prefers the destructive outcome  $(H, H)$  over the remaining outcomes.

Under the assumption that  $g > l$  and that the marginal cognitive costs are sufficiently small (but non-vanishing), we construct an ESC in which only individuals with preferences from  $\{u^i\}_{i=1}^{\infty}$  are present. Individuals of different cognitive levels coexist, and non-Nash profiles are played in all matches between equals. When individuals of the same level meet, they play mutual cooperation  $(D, D)$ . When individuals of different levels meet, the higher level plays  $H$  and the lower level plays  $D$ . The gain from obtaining the high payoff of  $1 + g$  against lower types is exactly counterbalanced by the higher cognitive costs. In contrast, if  $g < l$  then the game does not admit this kind of stable configuration.

**Proposition 10.** *Let  $G$  be the game represented in (2), where  $g > 0$  and  $l \in (0, 1)$ . Suppose that the marginal cognitive cost is small but non-vanishing, so that (a) there is an  $N$  such that  $k_N \leq l + g < k_{N+1}$ , and (b) it holds that  $g > k_{n+1} - k_n$  for all  $n \leq N$ .*

(i) *If  $g > l$  then there exists an ESC  $(\mu^*, b^*)$ , such that  $C(\mu^*) \subseteq \{(u^n, n)\}_{n=1}^N$ , and  $\mu^*$  is mixed (i.e.,  $|C(\mu^*)| > 1$ ). The behaviour of the incumbents is as follows: if the individuals in a match are of different cognitive levels, then the higher level plays  $H$  and the lower level plays  $D$ ; if both individuals in a match are of the same cognitive level, then they both play  $D$ .*

(ii) *If  $g = l$  then there exists an NSC with the above properties.*

(iii) *If  $g < l$  then there does not exist any NSC  $(\mu^*, b^*)$ , such that  $C(\mu^*) \subseteq \{(u^n, n)\}_{n=1}^{\infty}$ .*

The formal proof is presented in Appendix C.3.

## 7 Conclusion and Directions for Future Research

We have developed a model in which preferences coevolve with the ability to detect others' preferences and misrepresent one's own preferences. To do so, we have allowed for heterogeneity with respect to costly cognitive ability. The assumption of an exogenously given level of observability of the opponent's preferences, which has characterized the indirect evolutionary approach up until now, is replaced by a Machiavellian notion of deception equilibrium, which endogenously determines what each player observes. Only when players in a match are of the same cognitive level do they mutually observe each other's preferences. Our main results surprisingly show that this "grain" of perfect observability among equals is enough to imply that only efficient configurations can be stable. Moreover, we show that if the marginal cognitive cost is not too high then only Nash equilibria can be stable. Previous research has shown that if preferences are perfectly (or almost perfectly) observable, then efficiency is necessary for stability (Dekel, Ely, and Yilankaya, 2007). Previous research has also shown that if preferences are not at all observable, then Nash equilibrium is necessary for stability (Ok and Vega-Redondo, 2001; Ely and Yilankaya, 2001). *In our model, with endogenous observability of preferences we find that both efficiency and Nash equilibrium are necessary conditions for stability (provided that the marginal cognitive cost is not too high).*

We verify that if players with equal levels cannot observe each other’s preferences, then stable configurations do not have to be efficient. Instead, they are closely related to Nash equilibria of the underlying games. This should not come as a surprise since this alternative tie-breaking rule removes the last grain of mutual observation.

Our model assumes a very powerful form of deception. This allows us to derive sharp results that clearly demonstrate effects of endogenising observation, and introducing deception. We think that the “Bayesian” deception is an interesting model for future research: each incumbent type is associated with a signal, agents with high cognitive levels can mimic the signals of types with lower cognitive levels, and agents maximize their preferences given the received signals and the correct Bayesian inference about the opponent’s type.

In a companion paper (Heller and Mohlin, 2015) we study environments in which players are randomly matched, and make inferences about the opponent’s type by observing her past behaviour (rather than observing the type directly as is standard in the “indirect evolutionary approach”). In future research, it would be interesting to combine both approaches and allow the observation of the past behaviour to be influenced by deception.

Most papers taking the indirect evolutionary approach study the stability of preferences defined over material outcomes. Moreover, it is common to restrict attention to some parameterised class of such preferences. Since we study preferences defined on the more abstract level of action profiles (or the joint set of action profiles and opponent’s types in the case of type-interdependent preferences) we do not make predictions about whether some particular kind of preferences over material outcomes, from a particular family of utility functions, will be stable or not. It would be interesting to extend our model to such classes of preferences. Furthermore, with preferences defined over material outcomes it would be possible to study coevolution of preferences and deception not only in isolated games, but also when individuals play many different games using the same preferences. We hope to come back to these questions and we invite others to employ and modify our framework in these directions.

## A Formal Proofs of Theorems 1–2

### A.1 Preliminaries

This subsection contains notation and definitions that will be used in the following proofs.

A generous action is an action such that if played by the opponent, it allows a player to achieve the maximal fitness payoff. Formally:

**Definition 16.** Action  $a_g \in A$  is *generous*, if there exists an  $a \in A$  such that  $\pi(a, a_g) \geq \pi(a', a'')$  for all  $a', a'' \in A$ .

Fix a generous action  $a_g \in A$  of the game  $G$ . A second-best generous action is an action such that if played by the opponent, it allows a player to achieve the fitness payoff which is maximal under the constraint that the opponent is not allowed to play the generous action  $a_g$ . Formally:

**Definition 17.** Action  $a_{g_2} \in A$  is *second-best generous*, conditional on  $a_g \in A$  being first-best generous, if there exists  $a \in A$  such that  $\pi(a, a_{g_2}) \geq \pi(a', a'')$  for all  $a', a'' \in A$  such that  $a'' \neq a_g$ .

Fix a generous action  $a_g \in A$ , and fix a second-best generous action  $a_{g_2} \in A$ , conditional on  $a_g \in A$  being



first-best generous. For each  $\alpha \geq \beta \geq 0$ , let  $u_{\alpha,\beta}$  be the following utility function:

$$u_{\alpha,\beta}(a, a') = \begin{cases} \alpha & a' = a_g \\ \beta & a' = a_{g_2} \\ 0 & \text{otherwise.} \end{cases}$$

Observe that such a utility function  $u_{\alpha,\beta}$  satisfies:

1. *Indifference*: the utility function only depends on the opponent's action; i.e., the player is indifferent between any two of her own actions.
2. *Pro-generosity*: the utility is highest if the opponent plays the generous action, second-highest if the opponent plays the second-best generous action, and lowest otherwise.

Let  $U_{GI} = \{u_{\alpha,\beta} | \alpha \geq \beta \geq 0\}$  be the family of all such preferences, called *pro-generous indifferent preferences*. Note that  $U_g$  includes a continuum of different utilities (under the assumption that  $G$  includes at least three actions). Thus for any set of incumbent types we can always find a utility function in  $U_g$  which does not belong to any of the current incumbents.

## A.2 Proof of Theorem 1 (Behaviour of the Highest Types)

### A.2.1 Proof of Theorem 1, Part 1

Assume to the contrary that  $\pi(b_{\bar{\theta}}(\bar{\theta}), b_{\bar{\theta}}(\bar{\theta})) < \bar{\pi}$ . (Note that the definition of  $\bar{\pi}$  implies that the opposite inequality is impossible.) Let  $a_1, a_2 \in A$  be any two actions such that  $(a_1, a_2)$  is an efficient action profile, i.e.,  $0.5 \cdot (\pi(a_1, a_2) + \pi(a_2, a_1)) = \bar{\pi}$ . Let  $\theta_1, \theta_2$  be two types that satisfy the following conditions: (1) the types are not incumbents:  $\theta_1, \theta_2 \notin C(\mu^*)$ , (2) both types have the highest incumbent cognitive level:  $n_{\theta_1} = n_{\theta_2} = \bar{n}$ , and (3) both types have different pro-generosity indifferent preferences:  $u_{\theta_1}, u_{\theta_2} \in U_{GI}$  and  $u_{\theta_1} \neq u_{\theta_2}$ . Let  $\mu'$  be the distribution that assigns mass 0.5 to each of these types. The post-entry type distribution is  $\tilde{\mu} = (1 - \epsilon) \cdot \mu + \epsilon \cdot \mu'$ . Let the post-entry behaviour policy  $\tilde{b}$  be defined as follows:

1. Behaviour among incumbents respects focality:  $\tilde{b}_{\theta}(\theta') = b_{\theta}(\theta')$  for each incumbent pair  $\theta, \theta' \in C(\mu^*)$ .
2. In matches between mutants and incumbents of lower types, behaviour is such that the mutants maximize their fitness:  $(\tilde{b}_{\theta_i}(\theta'), \tilde{b}_{\theta'}(\theta_i)) \in FMDE(\theta_i, \theta')$  for each  $i \in \{1, 2\}$  and  $\theta' \in C(\mu^*)$  with  $n_{\theta'} < \bar{n}$ . Note that  $FMDE(\theta_i, \theta')$  is nonempty in virtue of the construction of  $U_{GI}$ .
3. In matches between mutants and incumbents of the highest type, the mutants mimic  $\bar{\theta}$  and the incumbents of the highest type play the same way as they play against  $\bar{\theta}$ :  $(\tilde{b}_{\theta_i}(\theta'), \tilde{b}_{\theta'}(\theta_i)) = (b_{\bar{\theta}}(\theta'), b_{\theta'}(\bar{\theta}))$ , for each  $i \in \{1, 2\}$  and  $\theta' \in C(\mu^*)$  with  $n_{\theta'} = \bar{n}$ .
4. Two mutants of *different* types play efficiently when meeting each other:  $\tilde{b}_{\theta_i}(\theta_j) = a_i$  for each  $i \neq j \in \{1, 2\}$ .
5. Two mutants of the *same* type play like  $\bar{\theta}$  plays against itself, when meeting each other:  $\tilde{b}_{\theta_i}(\theta_i) = b_{\bar{\theta}}(\bar{\theta})$  for each  $i \in \{1, 2\}$ .

In virtue of point 1 the construction  $(\tilde{\mu}, \tilde{b})$  is a focal configuration (with respect to  $(\mu, b)$ ). By point 2 the mutants  $\theta_1$  and  $\theta_2$  earn weakly more than  $\bar{\theta}$  against lower types. By point 3 the mutants earn exactly the same

as  $\bar{\theta}$  against all the highest incumbent types (including  $\bar{\theta}$ ). By points 4 and 5 the mutants on average earn strictly more than  $\bar{\theta}$  against the mutants. In total average fitness earned by  $\theta_1$  and  $\theta_2$  is strictly higher than that of  $\bar{\theta}$ , against a population that follows  $(\tilde{\mu}, \tilde{b})$ . This implies that  $\mu'$  is a strictly better reply against  $\mu^*$  in the population game  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ . Thus,  $\mu^*$  is not a symmetric Nash equilibrium, and therefore it is not an NSS, in  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ , which implies that  $\mu^*$  is not an NSC.

*Remark 4.* In this proof the mutants only outperform the incumbents on average. This allows for the possibility that one mutant earns strictly less than  $\bar{\theta}$ , while the other mutant earns strictly more than  $\bar{\theta}$ . That is enough to prove the desired result. However, we could also have proved our result with a construction involving three different mutants  $\theta_1, \theta_2, \theta_3$ , each of whom earns strictly more than  $\bar{\theta}$  in the post-entry configuration. This would be achieved if the three mutant types were equally frequent, and any two mutants of the *same* type played the same way  $\bar{\theta}$  plays against itself (point 5 above), but played as follows when facing another mutant: (i) mutant  $\theta_1$  plays  $a_1$  against  $\theta_2$  and  $a_2$  against  $\theta_3$ , (ii) mutant  $\theta_2$  plays  $a_1$  against  $\theta_3$  and  $a_2$  against  $\theta_1$ , and (iii) mutant  $\theta_3$  plays  $a_1$  against  $\theta_1$  and  $a_2$  against  $\theta_2$ .

### A.2.2 Proof of Theorem 1, Part 2

Assume to the contrary that  $((b_{\bar{\theta}}(\underline{\theta}), b_{\bar{\theta}}(\bar{\theta}))) \notin FMDE(\bar{\theta}, \underline{\theta})$ . Let  $\hat{\theta}$  be a type that satisfies: (1) not being an incumbent:  $\hat{\theta} \notin C(\mu^*)$ , (2) having the highest incumbent cognitive level:  $n_{\hat{\theta}} = \bar{n}$ , and (3) having pro-generous indifferent preferences:  $u_{\hat{\theta}} \in U_{GI}$ . Let  $\mu'$  be the distribution that assigns mass one to type  $\hat{\theta}$ . The post-entry type distribution is  $\tilde{\mu} = (1 - \epsilon) \cdot \mu + \epsilon \cdot \mu'$ . Let the post-entry behaviour policy  $\tilde{b}$  be defined as follows:

1. Behaviour among incumbents respects focality:  $\tilde{b}_{\theta}(\theta') = b_{\theta}(\theta')$  for each  $\theta, \theta' \in C(\mu^*)$ .
2. In matches between mutants and incumbents of lower types, behaviour is such that the mutants maximize their fitness:  $(\tilde{b}_{\hat{\theta}}(\theta'), \tilde{b}_{\theta'}(\hat{\theta})) \in FMDE(\hat{\theta}, \theta')$  for each  $\theta' \in C(\mu^*)$  with  $n_{\theta'} < \bar{n}$ .
3. In matches between mutants and incumbents of the highest type, the mutant mimics  $\bar{\theta}$  and the higher types play the same way they play against  $\bar{\theta}$ :  $(\tilde{b}_{\hat{\theta}}(\theta'), \tilde{b}_{\theta'}(\hat{\theta})) = (b_{\bar{\theta}}(\theta'), b_{\theta'}(\bar{\theta}))$ , for each  $\theta' \in C(\mu^*)$  with  $n_{\theta'} = \bar{n}$ .
4. The mutant  $\hat{\theta}$  plays against itself the same way  $\bar{\theta}$  plays against itself:  $(\tilde{b}_{\hat{\theta}}(\hat{\theta}), \tilde{b}_{\hat{\theta}}(\hat{\theta})) = (\tilde{b}_{\bar{\theta}}(\bar{\theta}), \tilde{b}_{\bar{\theta}}(\bar{\theta}))$ .

Note that  $(\tilde{\mu}, \tilde{b})$  is a focal configuration (with respect to  $(\mu, b)$ ), and that  $\hat{\theta}$  obtains a strictly higher fitness than  $\bar{\theta}$  against a population that follows  $(\tilde{\mu}, \tilde{b})$ . This implies that  $\mu'$  is a strictly better reply against  $\mu^*$  in the population game  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ . Thus,  $\mu^*$  is not a symmetric Nash equilibrium, and therefore it is not an NSS, in  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ , which implies that  $\mu^*$  is not an NSC.

### A.2.3 Proof of Theorem 1, Part 3

Assume to the contrary that  $\pi(b_{\bar{\theta}}(\bar{\theta}), b_{\bar{\theta}}(\underline{\theta})) > \bar{\pi}$ , which immediately implies that  $\pi(b_{\bar{\theta}}(\underline{\theta}), b_{\bar{\theta}}(\bar{\theta})) < \bar{\pi}$ . Let  $\hat{\theta}$  be a type that satisfies: (1) not being an incumbent:  $\hat{\theta} \notin C(\mu^*)$ , (2) having the highest incumbent cognitive level:  $n_{\hat{\theta}} = \bar{n}$ , and (3) having pro-generous indifferent preferences:  $u_{\hat{\theta}} \in U_{GI}$ . Let  $\mu'$  be the distribution that assigns mass one to type  $\hat{\theta}$ . The post-entry type distribution is  $\tilde{\mu} = (1 - \epsilon) \cdot \mu + \epsilon \cdot \mu'$ . Let the post-entry behaviour policy  $\tilde{b}$  be defined as follows:

1. Behaviour among incumbents respects focality:  $\tilde{b}_{\theta}(\theta') = b_{\theta}(\theta')$  for each  $\theta, \theta' \in C(\mu^*)$ .

2. In matches between mutants and incumbents of lower types, behaviour is such that the mutants maximize their fitness:  $(\tilde{b}_{\hat{\theta}}(\theta'), \tilde{b}_{\theta'}(\hat{\theta})) \in FMDE(\hat{\theta}, \theta')$  for each  $\theta' \in C(\mu^*)$  with  $n_{\theta'} < \bar{n}$ .
3. In a match between a mutant  $\hat{\theta}$  and the incumbent  $\bar{\theta}$ , the mutant mimics  $\underline{\theta}$ , and the incumbent  $\bar{\theta}$  plays the same way it plays against  $\underline{\theta}$ :  $(\tilde{b}_{\hat{\theta}}(\bar{\theta}), \tilde{b}_{\bar{\theta}}(\hat{\theta})) = (b_{\underline{\theta}}(\bar{\theta}), b_{\bar{\theta}}(\underline{\theta}))$ .
4. The mutant  $\hat{\theta}$  plays against itself the same way  $\bar{\theta}$  plays against itself:  $(\tilde{b}_{\hat{\theta}}(\hat{\theta}), \tilde{b}_{\hat{\theta}}(\hat{\theta})) = (\tilde{b}_{\bar{\theta}}(\bar{\theta}), \tilde{b}_{\bar{\theta}}(\bar{\theta}))$ .
5. The mutant  $\hat{\theta}$  mimics  $\bar{\theta}$  against all other highest types, and these higher types play against  $\hat{\theta}$  in the same way as they play against  $\bar{\theta}$ :  $(\tilde{b}_{\hat{\theta}}(\theta'), \tilde{b}_{\theta'}(\hat{\theta})) = (b_{\bar{\theta}}(\theta'), b_{\theta'}(\bar{\theta}))$  for each  $\theta' \neq \bar{\theta}$  with  $n_{\theta'} = \bar{n}$ .

Note that  $(\tilde{\mu}, \tilde{b})$  is a focal configuration (with respect to  $(\mu, b)$ ). By point 2 the mutant  $\hat{\theta}$  earns weakly more than  $\bar{\theta}$  against lower types. By point 3 and Theorem 1.1, the mutants earn strictly more than  $\bar{\theta}$  against type  $\bar{\theta}$ . By points 3 and 4 and Theorem 1.1, the mutant earns strictly more than  $\bar{\theta}$  against the mutant. By point 5 the mutant  $\hat{\theta}$  earns the same as  $\bar{\theta}$  against all other types. In total the average fitness earned by  $\hat{\theta}$  is strictly higher than that of  $\bar{\theta}$ , against a population that follows  $(\tilde{\mu}, \tilde{b})$ . This implies that  $\mu'$  is a strictly better reply against  $\mu^*$  in the population game  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ . Thus,  $\mu^*$  is not a symmetric Nash equilibrium, and therefore it is not an NSS, in  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ , which implies that  $\mu^*$  is not an NSC.

### A.3 Proof of Case (1) in Theorem 2

In what follows we fill in the missing technical details for the part of the proof of Theorem 2 that concerns case (A). We begin by proving a lemma.

**Lemma 1.** *If  $(\sigma_1, \sigma_2) \in DE(\theta_1, \theta_2)$  then there exist actions  $a_1, a'_1 \in C(\sigma_1)$  and  $a_2, a'_2 \in C(\sigma_2)$  such that:  $(a_1, a_2) \in DE(\theta_1, \theta_2)$ , and  $(a'_1, a'_2) \in DE(\theta_1, \theta_2)$ , with  $\pi(a_1, a_2) \geq \pi(\sigma_1, \sigma_2)$ , and  $\pi(a'_1, a'_2) \leq \pi(\sigma_1, \sigma_2)$ .*

*Proof.* Note that for any mixed deception equilibrium  $(\sigma_1, \sigma_2)$  and any action  $a \in C(\sigma_2)$ , the profile  $(\sigma_1, a)$  is also a deception equilibrium (because otherwise the deceiver would not induce the deceived party to take a mixed action that puts positive weight on  $a$ ). It follows that there are actions  $a_2, a'_2 \in C(\sigma_2)$  such that  $(\sigma_1, a_2)$  and  $(\sigma_1, a'_2)$  are deception equilibria, with  $\pi(\sigma_1, a_2) \geq \pi(\sigma_1, \sigma_2)$  and  $\pi(\sigma_1, a'_2) \leq \pi(\sigma_1, \sigma_2)$ . Furthermore, if  $(\sigma_1, a_2)$  and  $(\sigma_1, a'_2)$  are deception equilibria, then for any action  $a \in C(\sigma_1)$ , the profiles  $(a, a_2)$  and  $(a, a'_2)$  are also deception equilibria, with  $\pi(\sigma_1, a_2) = \pi(a, a_2)$  and  $\pi(\sigma_1, a'_2) = \pi(a, a'_2)$ . Hence there are actions  $a_1, a'_1 \in C(\sigma_1)$  such that  $(a_1, a_2)$  and  $(a'_1, a'_2)$  are deception equilibria, with  $\pi(a_1, a_2) = \pi(\sigma_1, a_2) \geq \pi(\sigma_1, \sigma_2)$ , and  $\pi(a_1, a'_2) = \pi(\sigma_1, a'_2) \leq \pi(\sigma_1, \sigma_2)$ .  $\square$

Assume that case (A) holds: there is an incumbent  $\hat{\theta}$  that plays inefficiently against itself, i.e.,  $(b_{\hat{\theta}}(\hat{\theta}), b_{\hat{\theta}}(\hat{\theta})) \neq (\bar{a}, \bar{a})$ , and there is no incumbent type with strictly higher cognitive level than  $\hat{\theta}$  that satisfies any of the cases (A), (B), or (C). To prove that this cannot hold in an NSC we introduce a mutant  $\hat{\theta} = (\hat{u}, n_{\hat{\theta}}) \notin C(\mu^*)$ . If  $\Sigma(u_{\hat{\theta}}) = \Delta$ , then we let  $\hat{u} \in U_{GI}$  be such that  $\hat{\theta} = (\hat{u}, n_{\hat{\theta}}) \notin C(\mu^*)$ . If  $\Sigma(u_{\hat{\theta}}) \neq \Delta$ , then we fix a dominated action  $\underline{a} \in A \setminus \Sigma(u_{\hat{\theta}})$ , and let  $\hat{u}$  be defined as follows:

$$\hat{u}(a, a') = \begin{cases} \max_{a \in A} (u_{\hat{\theta}}(a, \bar{a})) & a = a' = \bar{a} \\ u_{\hat{\theta}}(\underline{a}, a') - \beta_{a'} & a = \underline{a} \text{ and } a' \neq \bar{a} \\ u_{\hat{\theta}}(a, a') & \text{otherwise,} \end{cases}$$

where each  $\beta_{a'} \geq 0$  is chosen such that  $\hat{\theta} = (\hat{u}, n_{\hat{\theta}}) \notin C(\mu^*)$ . That is, if  $\Sigma(u_{\hat{\theta}}) \neq \Delta$ , then the utility function  $\hat{u}$  is constructed from the utility function  $u_{\hat{\theta}}$  by arbitrarily lowering the payoff of some of the outcomes associated with the (already) dominated action  $\underline{a}$  and which do not involve action  $\bar{a}$ , while increasing the payoff of the outcome  $(\bar{a}, \bar{a})$  by the minimal amount that makes  $\bar{a}$  a best reply to itself. Note that this definition of  $\hat{u}$  is valid also for the case of  $\bar{a} = \underline{a}$ . It follows that  $a \in \Sigma(u_{\hat{\theta}}) \cup \{\bar{a}\}$  iff  $a \in \Sigma(\hat{u})$ . To see this, note that if  $\Sigma(u_{\hat{\theta}}) \neq \Delta$  and  $\underline{a} = \bar{a}$ , then  $\Sigma(\hat{u}) = \Sigma(u_{\hat{\theta}}) \cup \{\bar{a}\}$ . Otherwise  $\Sigma(\hat{u}) = \Sigma(u_{\hat{\theta}})$ . Thus,  $\hat{\theta}$  can be induced to play exactly the same pure actions as  $\hat{\theta}$ , unless  $\bar{a} = \underline{a}$ , in which case  $\hat{\theta}$  can be induced to play  $\bar{a}$  in addition to all actions that  $\hat{\theta}$  can be induced to play.

Let  $\mu'$  be the distribution that assigns mass one to type  $(n_{\hat{\theta}}, \hat{u})$ . Let the post-entry type distribution be  $\tilde{\mu} = (1 - \epsilon) \cdot \mu + \epsilon \cdot \mu'$ , and let the post-entry behaviour policy  $\tilde{b}$  be defined as follows:

1. Behaviour among incumbents respects focality:  $\tilde{b}_{\theta}(\theta') = b_{\theta}(\theta')$  for each  $\theta, \theta' \in C(\mu^*)$ .
2. In matches between the mutant type  $\hat{\theta}$  and any lower type  $\theta' \in C(\mu^*)$  (with  $n_{\theta'} < n_{\hat{\theta}}$ ), we distinguish two cases.
  - (a) Suppose that  $\Sigma(u_{\hat{\theta}}) = \Delta$ . In this case let  $(\tilde{b}_{\hat{\theta}}(\theta'), \tilde{b}_{\theta'}(\hat{\theta})) \in FMDE(\hat{\theta}, \theta')$ . Note that  $FMDE(\hat{\theta}, \theta')$  is nonempty since in this case  $\hat{u} \in U_{GI}$ .
  - (b) Suppose that  $\Sigma(u_{\hat{\theta}}) \neq \Delta$ . In this case let  $(\tilde{b}_{\hat{\theta}}(\theta'), \tilde{b}_{\theta'}(\hat{\theta})) = (a_1, a_2)$ , for some  $(a_1, a_2) \in DE(\hat{\theta}, \theta')$  such that  $\pi(a_1, a_2) \geq \pi(b_{\hat{\theta}}(\theta'), b_{\theta'}(\hat{\theta}))$ . By Lemma 1 above such a profile  $(a_1, a_2)$  exists.
3. In matches between the mutant type  $\hat{\theta}$  and any incumbent type  $\theta'$  with same level, the mutant  $\hat{\theta}$  mimics  $\hat{\theta}$ , and the incumbent  $\theta'$  treats the mutant  $\hat{\theta}$  like the incumbent  $\hat{\theta}$ :  $(\tilde{b}_{\hat{\theta}}(\theta'), \tilde{b}_{\theta'}(\hat{\theta})) = (b_{\hat{\theta}}(\theta'), b_{\theta'}(\hat{\theta}))$  for all  $\theta'$  such that  $n_{\theta'} = n_{\hat{\theta}}$  and  $\theta' \neq \hat{\theta}$ .
4. The mutant plays efficiently when meeting itself:  $\tilde{b}_{\hat{\theta}}(\hat{\theta}) = \bar{a}$ .
5. In matches between the mutant type  $\hat{\theta}$  and a higher type  $\theta' \in C(\mu^*)$  (with  $n_{\theta'} > n_{\hat{\theta}}$ ), we distinguish two cases. Pick a profile  $(a_1, a_2) \in DE(\theta', \hat{\theta})$ , such that  $\pi(a_2, a_1) \geq \pi(b_{\hat{\theta}}(\theta'), b_{\theta'}(\hat{\theta}))$ . By Lemma 1 above such a profile  $(a_1, a_2)$  exists. Moreover, by the construction of  $\hat{u}$ , it is either the case that  $(a_1, a_2) \in DE(\theta', \hat{\theta})$ , or there is some  $\bar{a}$  such that  $u_{\theta'}(\bar{a}, \bar{a}) > u_{\theta'}(a_1, a_2)$ . In the latter case we have  $(\bar{a}, \bar{a}) \in DE(\theta', \hat{\theta})$ , due to the fact that  $(b_{\theta'}(\theta'), b_{\theta'}(\theta')) = (\bar{a}, \bar{a})$  implies that  $\bar{a}$  is a best reply to  $\bar{a}$  for type  $\theta'$ .
  - (a) If  $u_{\theta'}(a_1, a_2) > u_{\theta'}(\bar{a}, \bar{a})$  let  $(\tilde{b}_{\theta'}(\hat{\theta}), \tilde{b}_{\hat{\theta}}(\theta')) = (a_1, a_2)$ . Note that by definition of  $(a_1, a_2)$  it holds that  $\pi(a_2, a_1) \geq \pi(b_{\hat{\theta}}(\theta'), b_{\theta'}(\hat{\theta}))$ .
    - i. If  $u_{\theta'}(a_1, a_2) \leq u_{\theta'}(\bar{a}, \bar{a})$  let  $(\tilde{b}_{\theta'}(\hat{\theta}), \tilde{b}_{\hat{\theta}}(\theta')) = (\bar{a}, \bar{a})$ . Note that by definition of  $\hat{\theta}$  it holds that  $\pi(\bar{a}, \bar{a}) \geq \pi(b_{\hat{\theta}}(\theta'), b_{\theta'}(\hat{\theta}))$ .

By point 1,  $(\tilde{\mu}, \tilde{b})$  is a focal configuration (with respect to  $(\mu, b)$ ). By point 2 the mutant  $\hat{\theta}$  earns weakly more than  $\hat{\theta}$  against lower types. By point 3 the mutant  $\hat{\theta}$  earns the same as  $\hat{\theta}$  against all incumbents of level  $n_{\hat{\theta}}$ . By points 3 and 4 (and the assumption that  $\hat{\theta}$  does not play efficiently against itself), the mutant  $\hat{\theta}$  earns strictly more than  $\hat{\theta}$  against  $\hat{\theta}$ . By point 5 the mutant  $\hat{\theta}$  earns weakly more than  $\hat{\theta}$  against all incumbents of a higher cognitive level. In total the average fitness earned by  $\hat{\theta}$  is strictly higher than that of  $\hat{\theta}$ , against a population

that follows  $(\tilde{\mu}, \tilde{b})$ . This implies that  $\mu'$  is a strictly better reply against  $\mu^*$  in the population game  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ . Thus,  $\mu^*$  is not a symmetric Nash equilibrium, and therefore it is not an NSS of  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ , which implies that  $\mu^*$  is not an NSC. Thus we have shown that  $\theta^{\circ}$  plays efficiently against itself.

## B Variant with Uniform Deception

In this section we describe how to adapt our model in a way that requires players to use the *same* mixed action in their deception efforts towards all opponents with lower cognitive levels. We implement this change by replacing the definition of configuration with a new notion of *configuration with uniform deception*.

**Definition 18.** A *configuration with uniform deception* is a pair  $(\mu, b)$  where  $\mu \in \Delta(U)$  is a type distribution, and  $b : C(\mu) \times C(\mu) \rightarrow \Delta(A)$  is a behavioural policy such that

1. For each type  $\theta \in C(\mu)$ , there exists  $\tilde{\sigma}(\theta)$  that satisfies

$$\tilde{\sigma}(\theta) \in \arg \max_{\sigma \in \Delta(A)} \left( \sum_{\theta' \in C(\mu), n_{\theta'} < n_{\theta}} \mu(\theta') \cdot \max_{\sigma' \in BR_u(\sigma)} u_{\theta}(\sigma, \sigma') \right), \text{ and}$$

2. For each  $\theta, \theta' \in C(\mu)$ :  $b_{\theta}(\theta') = \tilde{\sigma}(\theta)$  and

$$n_{\theta} = n_{\theta'} \implies (b_{\theta}(\theta'), b'_{\theta'}(\theta)) \in NE(\theta, \theta'), \text{ and}$$

$$n_{\theta} > n_{\theta'} \implies b_{\theta'}(\theta) \in BR_{u_{\theta'}}(\tilde{\sigma}(\theta)).$$

We interpret  $\tilde{\sigma}(\theta)$  as the strategy that lower levels are deceived into believing is being played by type  $\theta$ , and we interpret  $b_{\theta}(\theta')$  as the strategy of type  $\theta$  when being matched with type  $\theta'$ .

We restrict our definition of a neutrally stable configuration to a configuration with uniform deceptions:

**Definition 19.** A configuration  $(\mu, b)$  is a *neutrally stable configuration (NSC) with uniform deception*, if for every  $\mu' \in \Delta(\Theta)$ , there is some  $\bar{\varepsilon} \in (0, 1)$  such that if  $(\tilde{\mu}, \tilde{b})$ , where  $\tilde{\mu} = (1 - \varepsilon) \cdot \mu + \varepsilon \cdot \mu'$ , is a focal configuration with uniform deceptions, then  $\mu$  is an NSS in the type game  $\Gamma_{(\tilde{\mu}, \tilde{b})}$ .

An analogous change can be made to the setup of interdependent preferences. All other details of the model are unchanged. It is relatively straightforward to see that *all our results hold also in this setup of uniform deceptions, with minor adaptations to the proofs*.

## C Constructions of Heterogeneous NSCs in Examples

The first subsection presents a lemma on stable heterogeneous populations, which will later be used to construct NSCs in the Rock-Paper-Scissors and Hawk-Dove games, with type-neutral and type-interdependent preferences.

## C.1 A Useful Lemma on Stable Heterogeneous Populations

Consider a configuration  $(\mu, b)$ , consisting of a type distribution with (finite) support  $C(\mu) \subseteq \{(u, n)\}_{n=1}^{\infty}$ , and behaviour policies such that

$$\pi(b_{\theta}(\theta'), b_{\theta'}(\theta)) = \begin{cases} t & \text{if } n_{\theta} > n_{\theta'} \\ w & \text{if } n_{\theta} = n_{\theta'} \\ s & \text{if } n_{\theta} < n_{\theta'} \end{cases} . \quad (3)$$

Thus  $t$  is the payoff that a player of type  $\theta$  earns when deceiving an opponent of type  $\theta'$ , and  $s$  is the payoff earned by the deceived party. When two individuals of the same type meet they earn  $w$ . Our first lemma concerns the type game  $\Gamma_{(\mu, b)}$  that is induced by a configuration  $(\mu, b)$ , such that  $C(\mu) \subseteq \{(u, n)\}_{n=1}^{\infty}$  and with behaviour policies given by (3). Although we have normalized  $k_1 = 0$  in the main text, we do not omit reference to  $k_1$  in what follows. This is done to simplify the proofs.

**Lemma 2.** *Suppose  $t \geq w \geq s$ . Suppose that there is an  $N$  such that*

$$k_N - k_1 \leq t - s < k_{N+1} - k_1, \quad (4)$$

and suppose that

$$t - w > k_{n+1} - k_n \text{ for all } n \leq N. \quad (5)$$

Consider the type game  $\Gamma_{(\mu, b)}$  induced by a configuration  $(\mu, b)$  with a type distribution such that  $C(\mu) \subseteq \{(u, n)\}_{n=1}^{\infty}$ , and with behaviour policies given by (3).

1. *If  $2w < s + t$  then  $\Gamma_{(\mu, b)}$  has a unique ESS  $\mu^* \in \Delta(C(\mu))$ , which is mixed, i.e.,  $C(\mu^*) > 1$ , and in which no type above  $N$  is present, i.e.,  $C(\mu^*) \subseteq \{(u, n)\}_{n=1}^N$ .*
2. *If  $2w = s + t$  then  $\Gamma_{(\mu, b)}$  has an NSS  $\mu^* \in \Delta(C(\mu))$ , which is mixed, i.e.,  $C(\mu^*) > 1$ , and in which no type above  $N$  is present, i.e.,  $C(\mu^*) \subseteq \{(u, n)\}_{n=1}^N$ .*
3. *If  $2w > s + t$  then  $\Gamma_{(\mu, b)}$ , admits no NSS and hence no ESS.*

The rest of this subsection is devoted to proving this result.

First note that type  $(u, N+1)$  earns strictly less than  $(u, 1)$  at all population states, and  $(u, N)$  earns at least as much as  $(u, 1)$  at least at some population state. This immediately follows from  $s \leq w \leq t$  and  $t - k_{N+1} < s - k_1$  and  $s - k_1 \leq t - k_N$ . For this reason it is sufficient to consider the type distributions with support in  $\{(u, n)\}_{n=1}^N$ . The payoffs for a type game with all these types present are

	$(u, 1)$	$(u, 2)$	$(u, 3)$	$\dots$	$(u, N-1)$	$(u, N)$	
$(u, 1)$	$w - k_1$	$s - k_1$	$s - k_1$	$\dots$	$s - k_1$	$s - k_1$	
$(u, 2)$	$t - k_2$	$w - k_2$	$s - k_2$	$\dots$	$s - k_2$	$s - k_2$	
$(u, 3)$	$t - k_3$	$t - k_3$	$w - k_3$	$\dots$	$s - k_3$	$s - k_3$	,
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$	
$(u, N-1)$	$t - k_{N-1}$	$t - k_{N-1}$	$t - k_{N-1}$	$\dots$	$w - k_{N-1}$	$s - k_{N-1}$	
$(u, N)$	$t - k_N$	$t - k_N$	$t - k_N$	$\dots$	$t - k_N$	$w - k_N$	

or in matrix form

$$\mathbf{A} = \begin{pmatrix} w - k_1 & s - k_1 & s - k_1 & \dots & s - k_1 & s - k_1 \\ t - k_2 & w - k_2 & s - k_2 & \dots & s - k_2 & s - k_2 \\ t - k_3 & t - k_3 & w - k_3 & \dots & s - k_3 & s - k_3 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ t - k_{N-1} & t - k_{N-1} & t - k_{N-1} & \dots & w - k_{N-1} & s - k_{N-1} \\ t - k_N & t - k_N & t - k_N & \dots & t - k_N & w - k_N \end{pmatrix}.$$

Inspecting the matrix  $\mathbf{A}$  we make the following observation:

*Claim 1.* Consider the game with payoff matrix  $\mathbf{A}$ . Suppose (5) holds.

1.  $(u, n + 1)$  is the unique best response to  $n$  for all  $n \in \{1, \dots, N - 2\}$ .

- (a) If  $t - k_N > s - k_1$  then  $(u, N)$  is the unique best reply to  $(u, N - 1)$ .
- (b) If  $t - k_N = s - k_1$  then  $(u, N)$  and  $(u, 1)$  are the only two best replies to  $(u, N - 1)$ .
- (c)  $(u, 1)$  is the unique best response to  $(u, N)$ .

*Proof.* Condition (5) implies that  $t - k_{N+1} > w - k_N$ , and the definition of  $N$  implies  $t - k_{N+1} < s - k_1$ . Taken together this implies that  $w - k_N < s - k_1$ , which means that  $(u, 1)$  is the unique best response to  $(u, N)$ .

The definition of  $N$  entails  $t - k_N \geq s - k_1$ . If  $t - k_N > s - k_1$  then  $(u, N)$  is the unique best reply to  $(u, N - 1)$ . If  $t - k_N = s - k_1$  then  $(u, N)$  and  $(u, 1)$  are the only two best replies to  $(u, N - 1)$ . Furthermore, (5) implies that  $(u, n + 1)$  is the unique best response to  $(u, n)$  for all  $n \in \{1, \dots, N - 2\}$ .  $\square$

It is an immediate consequence of the above lemma that all Nash equilibria of  $\mathbf{A}$  are mixed, i.e., that they have more than one type in their support. Next, we examine the stability properties of such equilibria. As discussed in the proof of Theorem 2, it is well-known that if  $\mathbf{A}$  is negative definite (semi-definite) with respect to the tangent space; i.e., if  $v \cdot \mathbf{A}v < 0$  for all  $v \in \mathbb{R}_0^d = \{v \in \mathbb{R}^d : \sum_{i=1}^d v_i = 0\}$ ,  $v \neq \mathbf{0}$ , then  $\mathbf{A}$  admits a unique ESS (but not necessarily a unique NSS). Moreover, the set of Nash equilibria coincides with the set of NSS and constitutes a nonempty convex subset of the simplex (Hofbauer and Sandholm 2009, Theorem 3.2).

One can show:

*Claim 2.* If  $2w \geq (\leq) s + t$  then  $\mathbf{A}$  is positive (negative) semi-definite w.r.t. the tangent space.

*Proof.* Let

$$\mathbf{K} = \begin{pmatrix} -k_1 & -k_1 & \dots & -k_1 \\ -k_2 & -k_2 & \dots & -k_2 \\ \vdots & \vdots & \ddots & \vdots \\ -k_N & -k_N & \dots & -k_N \end{pmatrix}, \quad \mathbf{B} = \begin{pmatrix} w & s & \dots & s \\ t & w & \dots & s \\ \vdots & \vdots & \ddots & \vdots \\ t & t & \dots & w \end{pmatrix},$$

so that

$$\mathbf{A} = \mathbf{B} + \mathbf{K}.$$

Note that  $v' \mathbf{K} v = 0$  for all  $v \in \mathbb{R}_0^N$ ,  $v \neq \mathbf{0}$ , so that  $v' \mathbf{A} v < 0$  for all  $v \in \mathbb{R}_0^N$ ,  $v \neq \mathbf{0}$ , if and only if  $v' \mathbf{B} v < 0$  for all  $v \in \mathbb{R}_0^N$ ,  $v \neq \mathbf{0}$ . Moreover, note that  $v' \mathbf{B} v < 0$  for all  $v \in \mathbb{R}_0^N$ ,  $v \neq \mathbf{0}$ , if and only if  $v' \bar{\mathbf{B}} v < 0$  for all  $v \in \mathbb{R}_0^N$ ,  $v \neq \mathbf{0}$ , where

$$\bar{\mathbf{B}} = \frac{1}{2} (\mathbf{B} + \mathbf{B}^T).$$

One can transform the problem to one of checking negative definiteness with respect to  $\mathbb{R}^{N-1}$  rather than the tangent space  $\mathbb{R}_0^N$ ; see, e.g., Weissing (1991). This is done with the  $N \times (N - 1)$  matrix  $\mathbf{P}$  defined by

$$p_{ij} = \begin{cases} 1 & \text{if } n = j \text{ and } n, j < N \\ 0 & \text{if } n \neq j \text{ and } n, j < N \\ -1 & \text{if } n = N \end{cases} .$$

We have

$$\mathbf{P}'\bar{\mathbf{B}}\mathbf{P} = \left( w - \frac{1}{2}(s + t) \right) (\mathbf{I} + \mathbf{1}\mathbf{1}'),$$

where  $\mathbf{1}$  is an  $N - 1$ -dimensional vector with all entries equal to 1, and  $I$  is the identity matrix. The matrix  $\mathbf{P}'\bar{\mathbf{B}}\mathbf{P}$  has one eigenvalue (of multiplicity  $N - 1$ ) that is equal to  $2w - (s + t)$ . Finally, note that this eigenvalue is non-negative if and only if  $2w \geq (s + t)$ .  $\square$

It follows that if  $2w \leq s + t$  then the game with payoff matrix  $\mathbf{A}$  admits an NSS. If  $2w > s + t$  then the game does not have a mixed NSS. We are now able to prove Lemma 2.

1. If  $2w < s + t$  then by Lemma 2  $\mathbf{A}$  is negative definite w.r.t. the tangent space, implying that it has a unique ESS. Lemma 1 implies that there can be no pure Nash equilibria (and hence no pure ESS). Thus  $\mathbf{A}$  has a unique Nash equilibrium, which is mixed. As observed earlier, type  $(u, N + 1)$  (and higher types) earn strictly less than  $(\theta, 1)$  for all population states, which implies that this unique equilibrium remains an ESS also when they are included in the set of feasible types.
2. If  $2w = s + t$  then  $\mathbf{A}$  is both positive and negative semi-definite w.r.t. the tangent space. In this case  $\mathbf{A}$  does not have an ESS but it does have a set of NSSs, all of which are Nash equilibria. Moreover, we know that  $\mathbf{A}$  has no pure NE, and so all NSS are mixed. Again, type  $(u, N + 1)$  (and higher types) can be ignored because they always earn strictly less than  $(\theta, 1)$ .
3. If  $2w > s + t$  then  $\mathbf{A}$  is positive definite w.r.t. the tangent space, implying that it has no NSC.

## C.2 Proof of Proposition 4: Equilibrium in Rock-Paper-Scissors Game

Formally the behaviour of the incumbent types is as follows:

$$b_{\theta}^*(\theta') = \begin{cases} (0, 1, 0) & \text{if } n_{\theta} > n_{\theta'} \\ (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}) & \text{if } n_{\theta} = n_{\theta'} \\ (1, 0, 0) & \text{if } n_{\theta} < n_{\theta'} \end{cases} .$$

Under the described behavioural policy we have

$$\pi(b_{\theta}(\theta'), b_{\theta'}(\theta)) = \begin{cases} 1 & \text{if } n_{\theta} > n_{\theta'} \\ 0 & \text{if } n_{\theta} = n_{\theta'} \\ -1 & \text{if } n_{\theta} < n_{\theta'} \end{cases} .$$

Start by restricting attention to the set of types  $\{(u^{\pi}, n)\}_{n=1}^{\infty}$ . That is, for the moment we use  $\{(u^{\pi}, n)\}_{n=1}^{\infty}$  instead of  $\Theta$  as the set of all types. All definitions can be amended accordingly. Lemma 2 in Appendix B implies that there is an NSC  $(\mu^*, b^*)$ , such that  $C(\mu^*) \subseteq \{(u^{\pi}, n)\}_{n=1}^N$ , and  $\mu^*$  is mixed. Lemma 2 establishes that the type game between the types  $\{(u^{\pi}, n)\}_{n=1}^N$  behaves much like an  $N$ -player version of a Hawk-Dove game:



it has a unique symmetric equilibrium that is in mixed strategies and that is neutrally or evolutionarily stable, depending on whether the payoff matrix of the type game is negative semi-definite, or negative definite, with respect to the tangent space.

It remains to show that types not in  $\{(u^\pi, n)\}_{n=1}^\infty$  are unable to invade. Suppose a mutant of type  $(u', n')$  enters. Incumbents of level  $n > n'$  will give the mutant a belief that induces the mutant to play some action  $a'$  and then play action  $a' + 1 \pmod 3$ , which is the incumbents' best response to  $a'$ . Thus, against incumbents of level  $n > n'$  the mutant earns  $-1$ . Against incumbents of level  $n < n'$ , the mutant will earn at most 1. Against incumbents of level  $n'$  the mutant earns at most 0. Against itself the mutant (or a group of mutants for that matter) will earn 0. Thus any mutant of level  $n'$  earns weakly less than the incumbents of level  $n'$ , in any focal post-entry configuration.

*Remark 5.* Our analysis is similar to that of Conlisk (2001). Like us, he works with a hierarchy of cognitive types (though in his case it is fixed and finite), where higher cognitive types carry higher cognitive costs. He stipulates that when a high type meets a low type the high type gets 1 and the low type gets  $-1$ . If two equals meet both get 0. He shows that there is a neutrally stable equilibrium of this game between types (using somewhat different arguments than we do), and explores comparative static effects of changing costs. However, unlike in our model, in Conlisk's model all individuals have the same materialistic preferences and the payoffs earned from deception are not derived from an underlying game.

### C.3 Proof of Proposition 10: Equilibrium in Hawk-Dove with Type-Interdependent Preferences

Formally, the behaviour of an incumbent  $\theta \in C(\mu^*)$  facing another incumbent  $\theta' \in C(\mu^*)$  is given by

$$b_\theta^*(\theta') = \begin{cases} D & \text{if } n_\theta \geq n_{\theta'} \\ H & \text{if } n_\theta < n_{\theta'} \end{cases}. \quad (6)$$

Under the described behavioural policy we have, for  $\theta, \theta' \in \{(u^n, n)\}_{n=1}^\infty$ ,

$$\pi(b_\theta(\theta'), b_{\theta'}(\theta)) = \begin{cases} 1 + g & \text{if } n_\theta > n_{\theta'} \\ 1 & \text{if } n_\theta = n_{\theta'} \\ 1 - l & \text{if } n_\theta < n_{\theta'} \end{cases}.$$

Start by restricting attention to the set of types  $\{(u^n, n)\}_{n=1}^\infty$ . That is, for the moment, let  $\{(u^n, n)\}_{n=1}^\infty$ , instead of  $\Theta_{ID}$ , be the set of all types. All definitions can be amended accordingly. Under this restriction on the set of types, the desired results (i)–(iii) follow from Lemma 2 in Appendix B. For example, to see that Lemma 2 implies part (i) for the restricted type set, note that  $g > l$  implies that  $2w < t + s$ , and  $g > k_{n+1} - k_n$  implies that  $t - w > k_{n+1} - k_n$ , in the language of Lemma 2. The arguments for (ii) and (iii) are analogous.

Next, allow for a larger set of types  $\Theta_{ID}$ , such that  $\{(u^n, n)\}_{n=1}^\infty \subseteq \Theta_{ID}$ . The fact that part (iii) of Proposition 10 holds for the restricted set of types implies that it also holds for any larger set of types. It remains to prove parts (i) and (ii) for the full set of types. We prove only part (i). The proof of part (ii) is very similar.

Consider a population consisting exclusively of types from the set  $\{(u^n, n)\}_{n=1}^\infty$ , and assume that the type distribution of these incumbents, together with the behaviour policy (6), would have constituted an ESC if the type set had been restricted to  $\{(u^n, n)\}_{n=1}^\infty$ . Suppose a mutant of type  $(u', n') \notin \{(u^n, n)\}_{n=1}^\infty$  enters. If it is the case that type  $(u^{n'}, n')$  is not among the incumbents, then by the definition of an ESC, it must earn

weakly less against the incumbents than what the incumbents earn against each other. Thus it is sufficient to show that the mutant of type  $(u', n')$  earns less than what a mutant or incumbent of the same cognitive level, i.e., type  $(u^{n'}, n')$ , would earn.

Against an incumbent  $(u^n, n)$  of level  $n > n'$  a mutant of type  $(u', n')$  earns at most  $1 - l$ , and type  $(u^{n'}, n')$  earns  $1 - l$ . Against an incumbent  $(u^n, n)$  of level  $n = n'$  a mutant of type  $(u', n')$  earns at most  $1 - l$ , and type  $(u^{n'}, n')$  earns 1. Against incumbents  $(u^n, n)$  of level  $n < n'$  a mutant of type  $(u', n')$  earns at most  $1 + g$ , and type  $(u^{n'}, n')$  earns  $1 + g$ . Thus in all cases, a mutant  $(u', n') \notin \{(u^n, n)\}_{n=1}^{\infty}$  earns strictly less than what a mutant or incumbent of type  $(u^{n'}, n')$  earns. Hence if mutants are sufficiently rare they will earn strictly less than incumbents in any focal post-entry configuration.

## References

- ABREU, D., AND R. SETHI (2003): “Evolutionary Stability in a Reputational Model of Bargaining,” *Games and Economic Behavior*, 44(2), 195–216.
- ALGER, I., AND J. W. WEIBULL (2013): “Homo Moralis, Preference Evolution under Incomplete Information and Assortative Matching,” *Econometrica*, 81(6), 2269–2302.
- BANERJEE, A., AND J. W. WEIBULL (1995): “Evolutionary Selection and Rational Behavior,” in *Learning and Rationality in Economics*, ed. by A. Kirman, and M. Salmon, pp. 343–363. Blackwell, Oxford.
- BERGSTROM, T. C. (1995): “On the Evolution of Altruistic Ethical Rules for Siblings,” *American Economic Review*, 85(1), 58–81.
- BESTER, H., AND W. GÜTH (1998): “Is Altruism Evolutionarily Stable?,” *Journal of Economic Behavior and Organization*, 34, 193–209.
- BOLLE, F. (2000): “Is Altruism Evolutionarily Stable? And Envy and Malevolence? Remarks on Bester and Güth,” *Journal of Economic Behavior and Organization*, 42, 131–133.
- BOMZE, I. M., AND J. W. WEIBULL (1995): “Does Neutral Stability imply Lyapunov Stability?,” *Games and Economic Behavior*, 11(2), 173–192.
- CONLISK, J. (2001): “Costly Predation and the Distribution of Competence,” *American Economic Review*, 91(3), 475–484.
- CRESSMAN, R. (1997): “Local Stability of Smooth Selection Dynamics for Normal form Games,” *Mathematical Social Sciences*, 34(1), 1–19.
- DEKEL, E., J. C. ELY, AND O. YILANKAYA (2007): “Evolution of Preferences,” *Review of Economic Studies*, 74, 685–704.
- DUFWENBERG, M., AND W. GÜTH (1999): “Indirect Evolution vs. Strategic Delegation: A Comparison of Two Approaches to Explaining economic institutions,” *European Journal of Political Economy*, 15(2), 281–295.
- DUNBAR, R. I. M. (1998): “The Social Brain Hypothesis,” *Evolutionary Anthropology*, 6, 178–190.
- ELLINGSEN, T. (1997): “The Evolution of Bargaining Behavior,” *The Quarterly Journal of Economics*, 112(2), 581–602.

- ELY, J. C., AND O. YILANKAYA (2001): “Nash Equilibrium and the Evolution of Preferences,” *Journal of Economic Theory*, 97, 255–272.
- FERSHTMAN, C., AND Y. WEISS (1998): “Social Rewards, Externalities and Stable Preferences,” *Journal of Public Economics*, 70(1), 53–73.
- FRANK, R. H. (1987): “If Homo Economicus Could Choose his own Utility Function, Would He Want One with a Conscience?,” *The American Economic Review*, 77(4), 593–604.
- FRENKEL, S., Y. HELLER, AND R. TEPER (2014): “The endowment effect as a blessing,” mimeo.
- FRIEDMAN, D., AND N. SINGH (2009): “Equilibrium Vengeance,” *Games and Economic Behavior*, 66(2), 813–829.
- GAMBA, A. (2013): “Learning and Evolution of Altruistic Preferences in the Centipede Game,” *Journal of Economic Behavior and Organization*, 85(C), 112–117.
- GÜTH, W. (1995): “An Evolutionary Approach to Explaining Cooperative Behavior by Reciprocal Incentives,” *International Journal of Game Theory*, 24(4), 323–344.
- GÜTH, W., AND S. NAPEL (2006): “Inequality Aversion in a Variety of Games: An Indirect Evolutionary Analysis,” *The Economic Journal*, 116, 1037–1056.
- GÜTH, W., AND M. E. YAARI (1992): “Explaining Reciprocal Behavior in Simple Strategic Games: An Evolutionary Approach,” in *Explaining Process and Change*, ed. by U. Witt, pp. 22–34. University of Michigan Press, Ann Arbor, MI.
- GUTTMAN, J. M. (2003): “Repeated Interaction and the Evolution of Preferences for Reciprocity,” *The Economic Journal*, 113(489), 631–656.
- HEIFETZ, A., C. SHANNON, AND Y. SPIEGEL (2007): “What to Maximize if You Must,” *Journal of Economic Theory*, 133(1), 31–57.
- HELLER, Y. (2015): “Three steps ahead,” *Theoretical Economics*, 10, 203–241.
- HELLER, Y., AND E. MOHLIN (2015): “Observations on cooperation,” .
- HEROLD, F., AND C. KUZMICS (2009): “Evolutionary Stability of Discrimination under Observability,” *Games and Economic Behavior*, 67, 542–551.
- HINES, W. G. S., AND J. MAYNARD SMITH (1979): “Games Between Relatives,” *Journal of Theoretical Biology*, 79(1), 19–30.
- HOFBAUER, J. (2011): “Deterministic Evolutionary Game Dynamics,” in *Proceedings of Symposia in Applied Mathematics*, vol. 69.
- HOFBAUER, J., AND W. H. SANDHOLM (2009): “Stable Games and Their Dynamics,” *Journal of Economic Theory*, 144(4), 1665–1693.
- HOFBAUER, J., AND K. SIGMUND (1988): *The Theory of Evolution and Dynamical Systems*. Cambridge University Press, Cambridge.

- HOLLOWAY, R. (1996): "Evolution of the Human Brain," in *Handbook of Human Symbolic Evolution*, ed. by A. Lock, and C. R. Peters, pp. 74–125. Clarendon Press, Oxford.
- HOPKINS, E. (2014): "Competitive Altruism, Mentalizing and Signalling," *American Economic Journal: Microeconomics*, 6, 272–292.
- HUCK, S., AND J. OECHSSLER (1999): "The Indirect Evolutionary Approach to Explaining Fair Allocations," *Games and Economic Behavior*, 28, 13–24.
- HUMPHREY, N. K. (1976): "The Social Function of Intellect," in *Growing Points in Ethology*, ed. by P. P. G. Bateson, and R. A. Hinde, pp. 303–317. Cambridge University Press, Cambridge.
- KIM, Y.-G., AND J. SOBEL (1995): "An Evolutionary Approach to Pre-Play Communication," *Econometrica*, 63(5), 1181–1193.
- KIMBOROUGH, E. O., N. ROBALINO, AND A. J. ROBSON (2014): "The Evolution of "Theory of Mind": Theory and Experiments," Cowles Foundation Discussion Paper No. 1907R, Yale University.
- KINDERMAN, P., R. I. M. DUNBAR, AND R. P. BENTALL (1998): "Theory-of-Mind Deficits and Causal Attributions," *British Journal of Psychology*, 89, 191–204.
- KOÇKCESEN, L., AND E. A. OK (2000): "Evolution of Interdependent Preferences in Aggregative Games," *Games and Economic Behavior*, 31, 303–310.
- MATSUI, A. (1991): "Cheap-talk and Cooperation in a Society," *Journal of Economic Theory*, 54(2), 245–258.
- MAYNARD SMITH, J. (1982): *Evolution and the Theory of Games*. Cambridge University Press, Cambridge.
- MAYNARD SMITH, J., AND G. R. PRICE (1973): "The Logic of Animal Conflict," *Nature*, 246(5427), 15–18.
- MOHLIN, E. (2010): "Internalized Social Norms in Conflicts: An Evolutionary Approach," *Economics of Governance*, 11(2), 169–181.
- (2012): "Evolution of Theories of Mind," *Games and Economic Behavior*, 75(1), 299–312.
- NACHBAR, J. H. (1990): "Evolutionary Selection Dynamics in Games: Convergence and Limit Properties," *International Journal of Game Theory*, 19(1), 59–89.
- NORMAN, T. W. L. (2012): "Equilibrium Selection and the Dynamic Evolution of Preferences," *Games and Economic Behavior*, 74(1), 311–320.
- OK, E. A., AND F. VEGA-REDONDO (2001): "On the Evolution of Individualistic Preferences: An Incomplete Information Scenario," *Journal of Economic Theory*, 97, 231–254.
- POSSAJENNIKOV, A. (2000): "On the Evolutionary Stability of Altruistic and Spiteful Preferences," *Journal of Economic Behavior and Organization*, 42, 125–129.
- PREMACK, D., AND G. WOODRUFF (1979): "Does the Chimpanzee have a Theory of Mind," *Behavioral and Brain Sciences*, 1, 515–526.
- ROBSON, A. J. (1990): "Efficiency in Evolutionary Games: Darwin, Nash and the Secret Handshake," *Journal of Theoretical Biology*, 144(3), 379–396.

- ROBSON, A. J., AND L. SAMUELSON (2011): “The Evolutionary Foundations of Preferences,” in *The Social Economics Handbook*, ed. by J. Benhabib, A. Bisin, and M. Jackson, pp. 221–310. North Holland.
- RTISCHEV, D. (2012): “Evolution of Mindsight, Transparency and Rule-Rationality,” Unpublished.
- SAMUELSON, L. (1991): “Limit Evolutionarily Stable Strategies in Two-Player, Normal Form Games,” *Games and Economic Behavior*, 3(1), 110–128.
- (2001): “Introduction to the Evolution of Preferences,” *Journal of Economic Theory*, 97(2), 225–230.
- SANDHOLM, W. H. (2001): “Preference Evolution, Two-Speed Dynamics, and Rapid Social Change,” *Review of Economic Dynamics*, 4, 637–679.
- (2010): “Local Stability under Evolutionary Game Dynamics,” *Theoretical Economics*, 5(1), 27–50.
- SCHAFFER, M. E. (1988): “Evolutionarily Stable Strategies for a Finite Population and a Variable Contest Size,” *Journal of Theoretical Biology*, 132, 469–478.
- SCHELLING, T. C. (1960): *The Strategy of Conflict*. Harvard University Press.
- SCHLAG, K. H. (1993): “Cheap Talk and Evolutionary Dynamics,” Bonn Department of Economics Discussion Paper B-242.
- SELTEN, R. (1980): “A Note on Evolutionarily Stable Strategies in Asymmetric Animal Conflicts,” *Journal of Theoretical Biology*, 84(1), 93–101.
- SETHI, R., AND E. SOMANTHAN (2001): “Preference Evolution and Reciprocity,” *Journal of Economic Theory*, 97, 273–297.
- STAHL, D. O. (1993): “Evolution of Smart<sub>n</sub> Players,” *Games and Economic Behavior*, 5(4), 604–617.
- STENNEK, J. (2000): “The Survival Value of Assuming Others to be Rational,” *International Journal of Game Theory*, 29, 147–163.
- TAYLOR, P. D., AND L. B. JONKER (1978): “Evolutionary Stable Strategies and Game dynamics,” *Mathematical Biosciences*, 40(1–2), 145–156.
- THOMAS, B. (1985): “On Evolutionarily Stable Sets,” *Journal of Mathematical Biology*, 22(1), 105–115.
- WÄRNERYD, K. (1991): “Evolutionary Stability in Unanimity Games with Cheap Talk,” *Economics Letters*, 36(4), 375–378.
- (1998): “Communication, Complexity, and Evolutionary Stability,” *International Journal of Game Theory*, 27(4), 599–609.
- WEISSING, FRANZ, J. (1991): “Evolutionary Stability and Dynamic Stability in a Class of Evolutionary Normal Form Games,” in *Game Equilibrium Models I. Evolution and Game Dynamics*, ed. by R. Selten, pp. 29–97. Springer.