# A Folk Theorem with Codes of Conduct[*]

Juan I. Block[†]     David K. Levine[‡]

May 13, 2016

**Abstract**

We study self-referential games in which players can perfectly understand an opponent's strategy, which is referred to as a code of conduct. We show a folk theorem for finite normal form games. We also provide an application of codes of conduct: games that are played through computer programs.

## 1 Introduction

In this note we examine economic situations where players employ codes of conduct which are defined as a complete specification of how they play and their opponents should play. Players also receive signals about what code of conduct their opponent may be using, while their own code of conduct enables them to respond to these signals. We focus on the limit case of perfect revealing signals because we are interested in applications such as games played through agents whose codes of conduct are embedded in their compensation, if the agent is human, and in their programming, if automated.

We show a folk theorem for finite normal form games using simple trigger codes of conduct and under two observability assumptions. First, we assume that all players can observe their opponents' codes of conduct. Second, we examine the case in which a group of players observe rivals' codes of conduct but are allowed to send cheap talk messages in order to coordinate punishments. We conclude by discussing how codes of conduct can be applied to computer algorithms.

Conditional commitment devices may expand equilibrium payoffs (see e.g. (Tennenholtz, 2004), and (Kalai, Kalai, Lehrer, and Samet, 2010)). The ability to condition on opponents' devices allows for implementing any outcome where players receive at least their minmax payoffs. In contrast to these papers, we allow for mixed strategies and computer programs are a special case of the code of conduct space considered here. More recently, attention has been drawn to noisy environment; (Block and Levine, 2015) examine agents that observe imperfectly informative signals about each other codes of conduct, and (Bachi, Ghosh, and Neeman, 2014) study games in which deceptive agents may betray their true intentions.

[†]Faculty of Economics, University of Cambridge.

[‡]Departments of Economics, EUI and WUSTL.

## 2 The Model

We consider a finite *base game* $\Gamma = \{I, (S_i, u_i)_{i \in I}\}$ with set of players $I = \{1, \dots, N\}$. Each player $i$ chooses a strategy $s_i$ from the finite strategy set $S_i$. Let $s \in S := \times_i S_i$ be the corresponding profile of strategies. The preferences of player $i$ are represented by a utility function $u_i : S \to \mathbb{R}$.

For any base game $\Gamma$, we can define the *self-referential game* $G(\Gamma)$. Every player $i$ privately observes a signal $y_i \in Y_i$, where $Y_i$ is finite, and $y \in Y := \times_i Y_i$ is the set of private signal profiles. The strategy for player $i$ is referred to as a *code of conduct* denoted by $r^i$. Each player $i$ is endowed with a common space of codes of conduct $R_0 = \{r^i | r^i_j : Y_j \to S_j\}$, where each element of $r^i$ defines a mapping from $Y_j$ to $S_j$. We write $r \in R := \times_i R_0$ for the profile of codes of conduct.

For each profile of codes of conduct $r \in R$, let $\pi(\cdot | r)$ be the probability distribution over signal profiles $Y$. The collection of probability distributions over profile of private signals is given by $\{\pi(\cdot | r) \in \Delta(Y) | r \in R\}$. Let $\pi_i(\cdot | r)$ be the marginal probability distribution of $\pi(\cdot | r)$ over $Y_i$. That is, $\pi_i(y_i | r)$ is the probability that player $i$ observes $y_i \in Y_i$ if players have chosen profile of codes of conduct $r \in R$. The expected utility of player $i$ is $U_i(r) = \sum_{y \in Y} u_i(r^1_1(y_1), \dots, r^N_N(y_N)) \pi(y | r)$.

Before playing $\Gamma$ and observing any $Y_i$, players simultaneously choose $r^i$. Then each player $i$ privately observes $y_i$ and plays $r^i_i(y_i) = s_i$. A Nash equilibrium is a profile $r \in R$ such that for all players $i$ and any alternative code of conduct $\tilde{r}^i \neq r^i$, it follows that $U_i(r^i, r^{-i}) \geq U_i(\tilde{r}^i, r^{-i})$.

## 3 The Folk Theorem

We first assume that each player $i$ is able to detect all opponents that do not choose the same code of conduct. Specifically, we say that the self-referential game $G$ *permits detection* if for each player $i$ and all players $j \neq i$ there exists a set of private signals $\overline{Y}^i_j \subset Y_j$ such that for any profile of codes of conduct $r \in R$, and any $\tilde{r}^i \neq r^i$ we have $\pi_j(\overline{Y}^i_j | \tilde{r}^i, r^{-i}) = 1$ and $\pi_j(\overline{Y}^i_j | r) = 0$. We also define the minmax strategy against player $i$ by $\underline{s}^i_{-i} \in \arg\min_{s_{-i} \in S_{-i}} \max_{s_i \in S_i} u_i(s_i, s_{-i})$. Let $\underline{u}_i = u_i(\underline{s}^i_i, \underline{s}^i_{-i})$ where $\underline{s}^i_i$ denote $i$'s best response to $\underline{s}^i_{-i}$.

Our first result in the perfect information case is similar to (Levine and Pesendorfer, 2007) with the difference that we consider asymmetry and more than two players:

**Theorem 1.** *If $v_i = u_i(s) \geq \underline{u}_i$ for all players $i$ with strategy profile $s \in S$ and $G$ permits detection, then there exists an $r \in R$ such that $(v_1, \dots, v_N)$ is a Nash equilibrium payoff of $G(\Gamma)$.*

*Proof.* Take any profile $s \in S$ such that for any player $i$, $u_i(s) \geq \underline{u}_i$. Consider the code of conduct $\hat{r}^i \in R_0$ that prescribes

$$\hat{r}^i_j(y_j) := \begin{cases} s_j & \text{if } y_j \notin \overline{Y}^i_j, \\ \underline{s}^i_j & \text{otherwise.} \end{cases}$$

If all players choose this code of conduct, any player $i$ would get $U_i(\hat{r}) = u_i(s)$. Contrary, if player $i$ adheres to some $\tilde{r}^i$ so that $\tilde{r}^i_i(y_i) = \tilde{s}_i$ for all $y_i \in Y_i$ and any $\tilde{s}_j$; and $\tilde{r}^i_j(y_j) = s_j$ for all $y_j, s_j$, he gets $U_i(\tilde{r}^i, \hat{r}^{-i}) = \underline{u}_i$. It follows then that $\hat{r}$ is a Nash equilibrium of the self-referential game. $\quad\square$

We turn now to the case in which only some players receive these perfectly revealing signals. More precisely, the self-referential game is said to be *locally perfectly informative* if there is a subset of players $J \subset I$ and for all players $j \in J$ there exists a set of private signals $\overline{Y}_j^i \subset Y_j$ such that for any profile of codes of conduct $r \in R$, for each player $i$, and any $\tilde{r}^i \neq r^i$ we have $\pi_j(\overline{Y}_j^i | \tilde{r}^i, r^{-i}) = 1$ and $\pi_j(\overline{Y}_j^i | r) = 0$. All players $i \notin J$ observe a trivial signal $y^*$.

It is possible that the players who receive the perfectly revealing signal need to communicate to another player to implement a punish. We assume a cheap talk communication stage: After receiving private signals $y_i \in Y_i$, players send cheap talk signals, defined as *announcements*, $\tilde{y}_i \in \tilde{Y}_0$ with profile $\tilde{y} \in \tilde{Y} := \times_i \tilde{Y}_0$. Note that the identity of both the announcer and the recipient of that public message are crucial. Let $\tilde{y}_i^j \in \tilde{Y}_0$ be player $i$'s announcement pointing that opponent $j$ has deviated. We require that there be $N - 1$ of such possible announcements for each player $i$, that is, $\#\tilde{Y}_0 = N(N-1)$. We allow for not sending a message $\{\emptyset\} \in \tilde{Y}_0$. A *message* $m_i$ from player $i$ is a map chosen from $M_i = \{m_i | m_i : Y_i \to \tilde{Y}_0\}$ with message profile denoted by $m \in M := \times_i M_i$. Once announcements have been made, players choose a strategy $s_i \in S_i$ according to a map $\phi_i : \tilde{Y} \times Y_i \to S_i$ that belongs to the set of all such maps $\Phi_i$. A strategy for player $i$ in the base game $\Gamma$ with cheap talk is the decision about a message $m_i$ to send and a strategy $s_i \in S_i$ to take after observing $\tilde{Y}_0$, that is, $s_i' = (m_i, \phi_i) \in M_i \times \Phi_i$.

If player $j$ announces player $i$ plans to deviate from the code of conduct profile and this was the only announcement, all players might play the prescribed action required to implement punishment to player $i$. However, player $i$ may try to take advantage of this information structure by announcing somebody else has deviated. At worst when he is detected there will be two such announcements.[1] We rule out this mutually implication by assuming that the self-referential game *strongly permits detection*, meaning that for all pairs $i, j \in J$ there exists a set of private signals $\overline{Y}_j^i \subset Y_j$ such that for any profile of codes of conduct $r \in R$, and any $\tilde{r}^i \neq r^i$ we have $\pi_j(\overline{Y}_j^i | \tilde{r}^i, r^{-i}) = 1$ and $\pi_j(\overline{Y}_j^i) = 0$, but for any $\tilde{r}^j \neq r^j$ it holds that $\pi_i(y_i | \tilde{r}^j, r^{-j}) = \pi_i(y_i | r)$ for all $y_i \in Y_i$. What strong detection says in a sense is that there are "neutral" witnesses, that is, people who observe wrong-doing but who cannot be credibly accused of wrong-doing by the wrong-doer.

**Theorem 2.** *For all $s \in S$ such that $v_i = u_i(s) \geq \underline{u}_i$ for all $i$, if $G(\Gamma)$ strongly permits detection and is locally perfectly informative, then there is an $r \in R$ such that the payoff vector $v = (v_1, \ldots, v_N)$ is a Nash equilibrium of $G(\Gamma)$.*

*Proof.* Take $s \in S$ such that $u_i(s) \geq \underline{u}_i$ for all $i$. We construct $\hat{r}^i \in R_0$ by considering players $j \in J$ and those $j \notin J$ since $G$ is locally perfectly informative.

For $j \in J$, we define $r_j^i(y_j) = (m_j, \phi_j) \in M_j \times \Phi_j$ as follows. For all $y_j \in \overline{Y}_j^k$ with $j \neq k$, $m_j(y_j) = \tilde{y}_j^k$ for $\tilde{y}_j^k \in \tilde{Y}_0$ and for all $y_j \notin \overline{Y}_j^k$, $m_j(y_j) = \{\emptyset\}$. Next, for any $j \neq k$, $\phi_j(\tilde{y}, y_j) = \underline{s}_j^k$ if $\tilde{y}_j^k \in \tilde{y}$ and $y_j \in \overline{Y}_j^k$; and $\phi_j(\tilde{y}, y_j) = s_j$ if $\tilde{y}_j^k \notin \tilde{y}$ and $y_j \notin \overline{Y}_j^k$. For $j \notin J$, let $r_j^i(y_j) = (m_j, \phi_j) \in M_j \times \Phi_j$ be such that for all $j$, $m_j(y^*) = \{\emptyset\}$. For all $\tilde{y}_i^k \in \tilde{y}$ for some $i, k$, $\phi_j(\tilde{y}, y^*) = \underline{s}_j^k$ and for all $\tilde{y} = \{\emptyset\}$ it follows $\phi_j(\tilde{y}, y^*) = s_j$. It follows $U_i(\hat{r}) = u_i(s)$ for each $i$.

---

[1] This is a fairly common strategy in criminal proceedings: try to obscure guilt by blaming everyone else.

We begin by checking potential deviations of players $j \in J$, $\tilde{r}^j \neq \hat{r}^j$. It suffices to check the following cases. First, fix $\phi$ and pick any $m'_j \neq m_j$ such that $m'_j(y_j) = \{\emptyset\}$ for all $y_j \in Y_j$ implies that there will be only one $m_i(\overline{y}_i^j) = \tilde{y}_i^j \in \tilde{Y}_0$ for some $i \in J$ hence $U_j(\tilde{r}^j, r^{-j}) = \underline{u}_j$. Otherwise, take any $m'_j \neq m_j$ where $m'_j(y_j) = \tilde{y}_j^i \in \tilde{Y}_0$ for all $y_j \in Y_j$, it follows that there is also $m_i(y_i) = \tilde{y}_i^j \in \tilde{Y}_0$ for some $i \in J$, and by strongly permits detection player $j$ gets $U_j(\tilde{r}^j, r^{-j}) = \underline{u}_j$. Next, suppose that $m_j$ does not change but $\phi'_j \neq \phi_j$ such that $\phi'_j(\tilde{y}, y_j) = s'_j$ for $\tilde{y}_i^k \in \tilde{Y}_0$ for $i, k \neq j$ and $s'_j \neq s_j^k$ thereby implying there has been only one $\tilde{y}_i^j \in \tilde{Y}_0$ for $i \in J$ and $j$ obtains $U_j(\tilde{r}^j, r^{-j}) = \underline{u}_j$. Similarly, if $\phi'_j(\tilde{y}, y_j) = s'_j$ for any $\tilde{y} \in \tilde{Y}$ and $s'_j \neq s_j$, there is some $i \in J$ such that $m_i(y_i) = \tilde{y}_i^j \in \tilde{Y}_0$ and the punishment gives $U_j(\tilde{r}^j, r^{-j}) = \underline{u}_j$ to player $j$. We conclude by noting that for agents $j \notin J$, it must be that $m'_j = m_j$ since $Y_j = \{y^*\}$, while by $G$ being locally perfectly informative, either any $\phi'_j(\tilde{y}, y^*) = s'_j$ with $s'_j \neq s_j$ for $\tilde{y}_j^k \notin \tilde{y}$ or $\phi'_j(\tilde{y}, y^*) = s'_j$ with $s'_j \neq \underline{s}_j^k$ for $\tilde{y}_j^k \in \tilde{y}$, player $j$ receives $U_j(\tilde{r}^j, r^{-j}) = \underline{u}_j$. $\qquad \square$

# 4 Application: Codes of Conduct as Computer Algorithms

A physical model of strategies is to imagine that players play by submitting computer programs to play on their behalf. In the self-referential framework, computer programs work as follows. Fix a signal profile set $Y$ and break the program into two parts, one of which generates $Y$ based on analyzing the programs, the other of which maps the signal profiles $Y$ into the strategy profile set $S$. The programs are self-referential as they also receive as input the program of the other player.

Specifically, there is a finite language $L$ of computer statements, and a finite limit $l$ on the length of a program. The (finite) space of computer programs is $P = \{(x_n)_{n=1}^t \in L | t \leq l\}$. Each program $p^i \in P$ produces outputs which have the form of a map $p^i : P \times P \to \{1, 2, \ldots, \infty\} \times S$. The interpretation is that $p^i(p^1, p^2) = (\nu^i, s^i)$ produces the result $s^i$ after $\nu^i$ steps. In case $\nu^i = \infty$, the program does not halt. Notice that depending on the language $L$ these programs can be either Turing machines or finite state machines. A self-referential strategy is a pair consisting of a *default* strategy profile $\overline{s}^i \in S$ and a program, $r^i = (\overline{s}^i, p^i)$. After players submit their program $p^1, p^2$, each program $p^i$ is given itself and the program submitted by the opposing player $p^{-i}$ as inputs. All programs are halted after an upper limit of $\overline{\nu}$ steps. If $p^i(p^i, p^{-i}) = (\nu^i, s^i)$ and $\nu^i \leq \overline{\nu}$, that is, the program halted in time, we then define the mapping $r^i(p^1, p^2) = s^i$, otherwise $r^i(p^1, p^2) = \overline{s}^i$. To map this to a self-referential game we take the signal space to be $Y = S$. Finally, the probability distribution of signal profiles is $\pi(y|r) = 1$ if $y_i = r^i(p^i, p^{-i})$ for all players $i$, and $\pi(y|r) = 0$ otherwise.

**Definition** The strategy space $S$ is self-referential with respect to the deadline $\overline{\nu}$ if for every pair of actions $\overline{a}, \underline{a}$ there exists a strategy $s = (d, p)$ such that

$$p(\tilde{d}, \tilde{p}) := \begin{cases} \overline{\nu}, \overline{a} & \text{if } \tilde{d} = d, \tilde{p} = p, \\ \nu, \underline{a} & \text{otherwise.} \end{cases}$$

**Example 1.** We consider the trading game with common action space $A = \{0, 1\}$ and show that

the self-referential strategy space satisfies the properties of definition 4. The default action is $\overline{s}^i = 0$. The computer language is the Windows command language; the listing is given below:

```
@echo off
if "0" EQU "\%3" goto sameactions
echo 0
goto finish
:sameactions
echo n | comp \%2 \%4
if \%errorlevel\% EQU 0 goto cooperate
echo 0
goto finish
:cooperate
echo 1
:finish
```

This program runs from the Windows command line, and takes as inputs four arguments: a digit describing the own default action, a own filename, an opponent default action and an opponent filename. If the opponent default action is 0, and the opponent program is identical to the listing above, the program generates as its final output the number 1; otherwise it generates the number 0. The point is, since it has access to sequence of its own instructions, it simply compares them to the sequence of opponents program instructions to check if they are the same or not.

## 5    Conclusion

We showed a folk theorem for games where players observe perfectly informative signals that point at deviant codes of conduct and hence deviators are punished with certainty. We further weakened the assumption about who observe these signals, highlighting the importance of communication in self-referential games. We view codes of conduct in this specific environment as computer algorithms.

## References

Bachi, B., S. Ghosh, and Z. Neeman (2014). Communication and deception in 2-player games. *Working paper*.

Block, J. I. and D. K. Levine (2015). Codes of conduct, private information and repeated games. *International Journal of Game Theory*, forthcoming.

Kalai, A. T., E. Kalai, E. Lehrer, and D. Samet (2010). A commitment folk theorem. *Games and Economic Behavior 69*(1), 127–137.

Levine, D. K. and W. Pesendorfer (2007). The evolution of cooperation through imitation. *Games and Economic Behavior 58*(2), 293–315.

Tennenholtz, M. (2004). Program equilibrium. *Games and Economic Behavior 49*(2), 363–373.