

## Can A Turing Player Identify Itself?

David Levine  
*UCLA*

Balázs Szentes  
*University of Chicago*

### *Abstract*

We show that the problem of whether two Turing Machines are functionally equivalent is undecidable and explain why this is significant for the theory of repeated play and evolution.

---

We thank National Science Foundation Grants SES-03-14713 and SES-05-18762 for financial support.

**Citation:** Levine, David and Balázs Szentes, (2006) "Can A Turing Player Identify Itself?." *Economics Bulletin*, Vol. 1, No. 1 pp. 1-6

**Submitted:** January 29, 2006. **Accepted:** March 23, 2006.

**URL:** <http://www.economicsbulletin.com/2006/volume1/EB-06A00001A.pdf>

## 1. INTRODUCTION

The wide empirical prevalence of cooperation in prisoner's dilemma type situations where free-riding is possible has long fascinated economists, game theorists and biologists. The possibility of punishing free-riding in repeated situations is central to the theory of repeated games, but the folk-theorem leaves open the possibility of coordination failure and Pareto inferior equilibria. A key question addressed in the evolutionary literature is whether cooperation is merely possible, or whether in fact there are evolutionary forces that push in the direction of cooperation. There are a number of affirmative answers to the question: central to evolutionary forces that lead to cooperative outcomes is the ability to identify free-riders so they can be punished.<sup>1</sup>

There are three basic mechanisms through which "deviant" free-rider strategies can be detected: indirectly by observing behavior, through some sort of "cheap talk" in the form of a "secret handshake," or directly by observing the strategies themselves. The "secret handshake" is a signal that is secret from free-riders and is discussed by Robson [1990]. In the setting of play between finite state machines "secret handshakes" take place through observing play as described in Rubinstein [1986], for example. The possibility of directly observing strategies is examined in Levine and Pendorfer [2002] - they show a strong result that the only strategies to survive a stochastic process of evolution and imitation are strategies that cooperate when they meet themselves, and punish opponents who use different strategies.

Here we take up the issue of directly observing the strategies used by opponents. In other words, we envisage a setting where players exchange descriptions of their strategies prior to play. There are several reasons why this is interesting. First, whether information is generated through cheap-talk or through observation of behavior, it is not possible to have more information about an opponents intentions than through direct observation of her strategy. So direct observation gives an upper bound on how well a player can be informed about an opponent's intentions - any limitations when direct observation is used must also be limitations when intentions must be inferred indirectly. Second, there are interesting circumstances where direct observation of strategies is possible. For example, in a biological setting, strategies are encoded in DNA, and direct detection of DNA or its attributes is possible. A good example of this is the immune system which is based on detecting "free-riders" through their chemical attributes rather than their behavior. In the economic setting, institutions that make it easy to observe opponents strategic rules - for example, implementation by third parties, or the ability to credibly reveal the rules that are being used through institutions such as bureaucracies - are possible and indeed frequent.

Technically, a strategy is simply an algorithm that responds to certain inputs. Computer scientists conjecture that all algorithms that can be implemented and that meet certain minimal requirements are Turing machines,<sup>2</sup> so this is a natural way in which to discuss strategic algorithms.<sup>3</sup> Economists such as Binmore [1990] and Anderlini [1990] have explicitly modeled player strategies as Turing Machines: they have focused on showing that there does not exist a machine which

---

<sup>1</sup>In the literature on positive assortative matching, discussed, for example in Berstrom [2002], the matching mechanism implicitly identifies free-riders by matching players based on their play. This can be interpreted as indirectly detection through observed behavior.

<sup>2</sup>This is known as Church's Thesis.

<sup>3</sup>An alternative would be to consider finite state machines, also known as automata, a subclass of Turing Machines with finite memory. This would lead to a different set of issues: how much success could we expect a machine with  $n$  states to have in identifying the strategy of a machine with some larger number  $m > n$  states?

always gives a best-response to the opponent's strategy.<sup>4</sup> The issues can be illustrated through the example of matching pennies. Suppose that each player designs a Turing machine that takes as input the Turing machine designed by the other player. An obvious strategy would be to simulate the opponent's Turing machine, using your own Turing machine as input. The problem with this is that when both players use such Turing machines, we get an infinite regress in which I simulate you simulating me simulating you...and there is an infinite loop. In technical terms, these Turing machines do not halt when given each other as input, and so are useless in practice.

In the setting of punishing free-riders, the situation is less difficult. It is easy to design a machine that when given opposing Turing machine as input checks to see if that machine is the same or different.<sup>5</sup> The strategy is then to cooperate if the opposing machine is the same and punish if it is different. This is the basis for the Pendorfer-Levine [2002] result, for example. However, the same strategy may be described by many different Turing machines. For example, one machine might engage in some irrelevant preliminary calculations that are later discarded, while another does not. This raises the question whether it is possible to identify an opponent not by the details of how their strategy is described, but according to the strategy or function that their algorithm implements. That is the question addressed in this note. It is important for three reasons. First, the interesting question about opponent play is how they will behave, not how their code is written. The question we would like to know from examining an opponent's strategy are of the type: Is my opponent going to cheat, so I should do likewise? Is he going to engage in a strategy that leaves him open to exploitation? Second, there is an obvious social disadvantage if machines that differ merely in external appearance engage in costly punishment against one another. Third, if strategies can only be determined through observing behavior, then it is evident that only the function of a machine can be identified.

The questions we will examine, then, is the extent to which it is possible to identify an algorithm based on its function, rather than on its description.<sup>6</sup> We are not going to require Turing Machines to predict the strategies of their opponents, but focus on the narrower question at issue in sustaining cooperation: can Turing Machine recognize whether their opponents are computationally equivalent to themselves?

## 2. THE QUESTION

We let  $\tau$  denote a Turing Machine, and  $\tau(x)$  denote the output of  $\tau$  on input the integer  $x$ ,<sup>7</sup> where  $\tau(x)$  is an integer, or  $\tau(x) = \infty$  if  $\tau$  does not halt on input  $x$ . By convention, since there are countably many Turing Machines, we identify their descriptions with integers. We let  $\langle \tau \rangle$  denote the description of the Turing Machine  $\tau$ . In other words  $\tau'(\langle \tau \rangle)$  is the output of the Turing Machine  $\tau'$  when the input is the description of the Turing Machine  $\tau$ . If we have two Turing Machines  $\tau_1$  and  $\tau_2$  we say that they are *equivalent* if they compute the same function, that is  $\tau_1(x) = \tau_2(x)$  for all  $x$ . We say that a Turing Machine is *finite* if it halts on every input. A problem with a yes/no

---

<sup>4</sup>Canning [1992] shows that if the domain of possible players and games are appropriately restricted then Turing Machines will play Nash Equilibria. These restrictions are fairly severe.

<sup>5</sup>This is obvious for Turing Machines. Since a program that checks an input to see if it is identical to itself requires only enough memory to store itself and the input, it is also true for finite state machines.

<sup>6</sup>An alternative approach is to consider Bayesian priors as in Nachbar [1997].

<sup>7</sup>Occasionally we refer to Turing machines that take as input pairs of integers; since there is a computable 1-1 map from pairs of integers to integers, such a machine can be understood as taking as input the single integer generated by a fixed map of this type.

answer is *decidable* if there is a Turing machine that provides the output zero if the answer is no and the output one if the answer is yes.

We consider a series of questions about the decidability of the functional equivalence of two Turing machines. In each case, we ask if there is a Turing machine that can determine whether two Turing machines are computationally equivalent. The answer to each question is no: the problem is not decidable, or in other words, no Turing machine can determine equivalence in any of these cases.

**Question 1:** Is it a decidable problem whether two Turing Machines are equivalent?

**Question 2:** Is it a decidable problem whether two finite Turing Machines are equivalent?

**Question 3:** Is there some Turing Machine  $\tau$  for which it is a decidable problem whether a Turing Machine is equivalent to  $\tau$ ?

**Question 4:** Is there some finite Turing Machine  $\tau$  for which it is a decidable problem whether another finite Turing Machine is equivalent to  $\tau$ ?

We should emphasize that the problem of functional equivalence of Turing machines has not escaped the notice of computer scientists. For example, Question 1 is essentially a restatement of Rice's Theorem.<sup>8</sup> Since the proofs are short, and the terminology of "index sets" used in Rice's Theorem may not be familiar, we give simple direct proofs that the answers to all of these questions is "no."

### 3. THE ANSWER

Our point of departure is the famous Halting Lemma that determining whether or not a Turing Machine halts is not a decidable problem - that is, there is no Turing Machine which can take as input the description of a Turing Machine and an integer  $n$  and compute 1 if it halts on  $n$  and 0 if it does not. In other words, whether  $\tau(\langle\tau\rangle)$  halts or not is not decidable. For expositional purposes we provide a proof for this result from which the answers to Questions 1-4 will follow.

**Halting Lemma.** *Whether a Turing Machine halts on itself is undecidable.*

*Proof.* Suppose by contradiction, that there exists a Turing Machine  $\tau^*$  such that  $\tau^*(\langle\tau\rangle) = 1$  if  $\tau(\langle\tau\rangle)$  halts and zero otherwise. Construct the Turing Machine  $\tau'$  as follows.

$$\tau'(\langle\tau\rangle) = \begin{cases} 0 & \text{if } \tau^*(\langle\tau\rangle) = 0 \\ \infty & \text{if } \tau^*(\langle\tau\rangle) = 1. \end{cases}$$

Then  $\tau'(\langle\tau'\rangle) = 0 \Leftrightarrow \tau^*(\langle\tau'\rangle) = 0$  by the construction of  $\tau'$ . However, by the definition of  $\tau^*$ ,  $\tau^*(\langle\tau'\rangle) = 0$  is true if and only if  $\tau'$  does not halt, that is,  $\tau'(\langle\tau'\rangle) = \infty$ . So  $\tau'(\langle\tau'\rangle) = 0 \Leftrightarrow \tau'(\langle\tau'\rangle) = \infty$  an obvious contradiction.  $\square$

We now answer each question in the negative by showing that if it had a positive answer then we could find a Turing Machine that determines whether  $\tau(\langle\tau\rangle)$  halts for any Turing Machine  $\tau$ , contradicting the fact that this problem is undecidable.

In addition to the Halting Lemma, the other key ingredient in our proofs will be the idea of one Turing machine simulating another. That is, if we are given as input the description of a Turing machine  $\langle\tau\rangle$  and an input  $x$  we can as part of our design simply execute  $\tau(x)$  to see what the input Turing machine  $\tau$  would do on the input  $x$ . Of course,  $\tau(x)$  might not halt, in which case the machine we are designing would not halt either.

---

<sup>8</sup>See Corollary 6.1.5 in Cutland [1980].

**Answer to Question 1:** *Is it a decidable problem whether two Turing Machines are equivalent?* No. Fix a Turing Machine  $\tau$ . Define  $\tau'$  as follows. If  $x \neq \langle \tau \rangle$  then  $\tau'(x) = \tau(x)$ . If  $x = \langle \tau \rangle$  then  $\tau'(\langle \tau \rangle) = \infty$ . To construct  $\tau'$  is easy. It is almost identical to  $\tau$ , except that first it checks whether the input is  $\langle \tau \rangle$  or not. If it is, it runs forever, otherwise it simulates  $\tau$ . Notice, that in order to define  $\tau'$  one does not need to know whether  $\tau(\langle \tau \rangle)$  halts or not. Observe, that  $\tau$  and  $\tau'$  are equivalent if and only if  $\tau(\langle \tau \rangle)$  does not halt. Hence, if there was a Turing Machine which determines whether  $\tau$  and  $\tau'$  are equivalent, it could also determine whether  $\tau(\langle \tau \rangle)$  halts or not.

**Answer to Question 2:** *Is it a decidable problem whether two finite Turing Machines are equivalent?* No. Fix a Turing Machine  $\tau$ . Define the Turing Machine  $\hat{\tau}$  as follows:  $\hat{\tau}(n) = 1$  if  $\tau(\langle \tau \rangle)$  did not halt after  $n$  steps, and zero otherwise. Let  $\tilde{\tau}$  be a Turing Machine such that  $\tilde{\tau}(n) = 1$  for all  $n$ . Notice that  $\hat{\tau}$  and  $\tilde{\tau}$  are finite, and in addition they are equivalent if and only if  $\tau(\langle \tau \rangle)$  does not halt. Hence a machine that could determine if two finite machines are equivalent could determine if  $\tau(\langle \tau \rangle)$  halts.

**Answer to Question 3:** *Is there some Turing Machine  $\tau$  for which it is a decidable problem whether a Turing Machine is equivalent to  $\tau$ ?* No. Before we start to answer the question, note that there is a Turing Machine  $\tau_s$  such that if  $\tau$  halts on at least one input, then  $\tau_s(\langle \tau \rangle)$  is an input on which  $\tau$  halts. The way to construct  $\tau_s$  is the following. First  $\tau_s$  simulates the first step of  $\tau$  on 1. Then if  $\tau_s$  simulated the  $n$ th step of  $\tau$  on  $k$ , for  $n \neq 1$ , it will simulate the  $n - 1$ th step of  $\tau$  on  $k + 1$ , while for  $n = 1$  it will simulate the  $(k + 1)$ th step on input 1. In other words

step	input
1	1
2	1
1	2
3	1
2	2
1	3
4	1

and so forth. Once  $\tau_s$  finds a step where  $\tau$  halts, it stops and its output is the input of  $\tau$  on which it halted.

Let  $\tau$  be some arbitrary Turing Machine, and let  $\tau_\infty$  denote a Turing Machine which does not halt on any input. Suppose by contradiction that there exists a Turing Machine  $\hat{\tau}$  such that it is decidable whether or not  $\hat{\tau}$  is equivalent to another Turing Machine. This would mean that there exists a Turing Machine  $\tau^*$  such that  $\tau^*(\langle \tau' \rangle) = 1$  if  $\tau'$  is equivalent to  $\hat{\tau}$  and zero otherwise.

We now construct a Turing Machine  $\tilde{\tau}$  as follows. The construction depends on which of the following two cases holds.

*Case 1:*  $\tau^*(\langle \tau_\infty \rangle) = 1$ . Then define  $\tilde{\tau}$  as follows. If  $x \neq 1$ ,  $\tilde{\tau}(x) = \hat{\tau}(x)$ . If  $x = 1$ , then  $\tilde{\tau}$  simulates  $\tau$  on  $\langle \tau \rangle$ . Thus  $\tilde{\tau}(1)$  halts if and only if  $\tau(\langle \tau \rangle)$  does.

*Case 2:*  $\tau^*(\langle \tau_\infty \rangle) = 0$ . Since this means that  $\hat{\tau}$  is not equivalent to  $\tau_\infty$ , it implies that  $\hat{\tau}$  halts for some input. Hence  $\tau_s(\langle \hat{\tau} \rangle)$  is well-defined. This enables us to define  $\tilde{\tau}(x) = \hat{\tau}(x)$  whenever  $x \neq \tau_s(\langle \hat{\tau} \rangle)$ . If  $x = \tau_s(\langle \hat{\tau} \rangle)$ , then

- *First:*  $\tilde{\tau}$  simulates  $\tau$  on  $\langle \tau \rangle$ . If it halts then
- *Second:*  $\tilde{\tau}$  simulates  $\hat{\tau}$  on  $x = \tau_s(\langle \hat{\tau} \rangle)$  and provides the same output as  $\hat{\tau}$ .

So  $\tilde{\tau}(\tau_s(\langle \hat{\tau} \rangle))$  halts if and only if  $\tilde{\tau}$  is functionally equivalent to  $\hat{\tau}$ .

We conclude that  $\tau(\langle \tau \rangle)$  halts if and only if  $(\tau^*(\tau_\infty), \tau^*(\tilde{\tau})) \in \{(1, 0), (0, 1)\}$ . So  $\tau^*$  can determine if  $\tau(\langle \tau \rangle)$  halts.

**Answer to Question 4:** *Is there some finite Turing Machine  $\tau$  for which it is a decidable problem whether another finite Turing Machine is equivalent to  $\tau$ ?* No. Suppose by contradiction that there exists a finite Turing Machine  $\hat{\tau}$  such that it is decidable whether or not  $\hat{\tau}$  is equivalent to another finite Turing Machine. Fix a Turing Machine  $\tau$  and construct a Turing Machine  $\tilde{\tau}$  as follows. On input  $n$ ,  $\tilde{\tau}$  simulates the first  $n$  steps of  $\tau$  on  $\langle \tau \rangle$ . Then  $\tilde{\tau}$  computes  $\hat{\tau}(n)$ . If  $\tau$  halts before the  $n$ th step on  $\langle \tau \rangle$  then  $\tilde{\tau}(n) \neq \hat{\tau}(n)$ , otherwise  $\tilde{\tau}(n) = \hat{\tau}(n)$ . Obviously,  $\tilde{\tau}$  is finite because  $\hat{\tau}$  is also finite. Furthermore,  $\tilde{\tau}$  and  $\hat{\tau}$  are equivalent if and only if  $\tau(\langle \tau \rangle)$  does not halt.

#### REFERENCES

- [1] Anderlini, L. [1990], "Some Notes on Church's Thesis and the Theory of Games," *Theory and Decision* **29**, 19-52.
- [2] Bergstrom, T. [2002], "Evolution of Social Behavior: Individual and Group Selection," *The Journal of Economic Perspectives* **16**, 67-88.
- [3] Binmore, K. [1990], *Essays on the Foundation of Game Theory*. Oxford: Basil Blackwell.
- [4] Canning, D. [1992], "Rationality, Computability, and Nash Equilibrium," *Econometrica* **60**, 877-888.
- [5] Cutland, N. J. [1980], *Computability: An introduction to Recursive Function Theory*, Cambridge University Press.
- [6] Levine, D. K. and W. Pesendorfer [2002], "Evolution of Cooperation Through Imitation," UCLA.
- [7] Nachbar, J. [1997], "Prediction, Optimization, and Learning in Repeated Games," *Econometrica* **65**, 275-309.
- [8] Robson, A. J. [1990], "Efficiency in evolutionary games: Darwin, Nash and the secret handshake," *Journal of Theoretical Biology* **144**, 379-96.
- [9] Rubinstein A. [1986], "Finite Automata Play the Repeated Prisoner's Dilemma," *Journal of Economic Theory* **39**, 83-96.