

# NUMERICAL SIMULATIONS OF REINFORCEMENT WITH ITERATIVE PUNISHMENT

Jeffrey A. Barrett

Department of Logic and Philosophy of Science, UC Irvine  
j.barrett@uci.edu

Nathan Gabriel

Department of Logic and Philosophy of Science, UC Irvine  
ngabrie2@uci.edu

**ABSTRACT.** We consider the efficacy of various forms of reinforcement learning with punishment in evolving linguistic conventions in the context of Lewis-Skyrms signaling games. We show that the learning strategy of reinforcement with iterative punishment is highly effective at evolving optimal conventions in even complex signaling games. It is also robust and can be easily extended to a self-tuning variety of reinforcement learning. We briefly discuss some of the virtues of reinforcement with iterative punishment and how it may be related to learning in nature.

## 1. INTRODUCTION

We are concerned here with reinforcement learning with punishment in the context of Lewis-Skyrms signaling games.<sup>1</sup> In such games, punishment often contributes to both the speed and accuracy of reinforcement learning.<sup>2</sup>

It has been shown that a signaling game with two states, two signals, and two acts will converge to optimal signaling on simple reinforcement learning with no punishment if nature is unbiased.<sup>3</sup> But if one has more than two states, signals, and acts or if nature is biased, then the players often get stuck in a suboptimal pooling equilibrium. This problem becomes more pronounced in more complex games with more states, signals, and acts.<sup>4</sup>

Supplementing reinforcement on success with punishment on failure often facilitates the evolution of optimal player dispositions, but there is a tuning problem. If one punishes too much relative to the level of reinforcement on success, it makes it difficult for the players to learn anything at all. And if one punishes too little, one does not do much to solve the suboptimal pooling equilibrium problem. What

---

<sup>1</sup>David Lewis (1969) first presented this sort of game in the context of classical game theory as a way of studying how conventions might be established. Skyrms (2006, 2010) translated Lewis' signaling games into the context of evolutionary game theory. See Herrnstein (1970), Roth and Erev (1995), and Erev and Roth (1998) for a description of reinforcement learning and how it maybe used to characterize human learning in a number of salient contexts.

<sup>2</sup>See Barrett and Zollman (2009) for a discussion of the role of punishment in reinforcement learning. Punishment here just means weakening of dispositions by means of negative reinforcement.

<sup>3</sup>See Argiento, Pemantle, Skyrms, and Volkov (2009) for this proof. The proof is nontrivial for even this simple case. There are very few analytic results for more complex signaling games, under even simple reinforcement, that they must be investigated primarily by simulation.

<sup>4</sup>See Barrett (2006) for an early description of both of these problems and Hofbauer and Huttegger (2008) for further discussions of the latter.

counts as too much and too little here depends on the complexity of the game. This makes it difficult to specify a single dynamics that does well in a broad assortment of games.<sup>5</sup>

We will begin by considering a  $4 \times 4 \times 4$  signaling game with basic reinforcement learning, then consider simple reinforcement with punishment, reinforcement with iterative punishment, and self-tuning reinforcement with iterative punishment.<sup>6</sup> We will show how the latter form of reinforcement with punishment provides robust learning strategies that are effective in very different signaling games. How the virtues of reinforcement with iterative punishment pave the way for a self-tuning form of reinforcement learning will be particularly salient. We conclude by briefly discussing how reinforcement with iterative punishment may be related to learning in natural contexts.

## 2. SIMPLE REINFORCEMENT

A Lewis-Skyrms signaling game is a common interest evolutionary game with one sender and one receiver. On a play of an  $n \times n \times n$  signaling game, the sender sees which of  $n$  states of nature obtains. She then sends one of  $n$  possible signals. The receiver cannot see the state of nature but can see the sender's signal. When she sees the signal on a play of the game, she performs one of  $n$  possible acts. Each act is appropriate to precisely one state of nature. If the receiver's act matches the current state of nature, then the play counts as a success; otherwise, it counts as a failure.

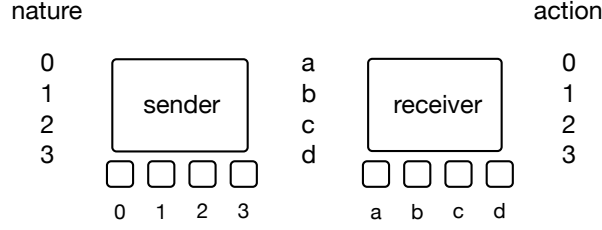
The two players update their conditional dispositions based on the outcome of each play (the sender updates her dispositions to signal conditional on the state and the receiver updates his dispositions to act conditional on the signal). Whether the agents are readily able to evolve optimal dispositions where each state of nature leads to a particular signal which leads to the act that matches the current state depends on precisely how they update their conditional dispositions given the outcome of each play. It sometimes often depends on luck as well.<sup>7</sup>

Let's consider concretely how this goes for a  $4 \times 4 \times 4$  signaling game under simple reinforcement learning (Figure 1). Simple reinforcement learning is characterized by stipulating how the probability of an agent's action on each round of a repeated interaction is updated given her successes and failures in earlier rounds. The model is being characterized by (1) an updating rule that determines how a agent's propensities evolve, (2) a response rule that determines how propensities determine dispositions, and (3) the initial propensities that determine a starting point for the process.

<sup>5</sup>See Barrett (2006) for an investigation of various punishment strategies and Barrett and Zollman (2009) for a discussion of the role of punishment and forgetting in learning. See Alexander, Skyrms, and Zabell (2012), Barrett, Cochran, Huttegger, and Fujiwara (2017), and Cochran and Barrett (2020) and (2021) for discussions of other learning strategies that seek to address the problem of suboptimal pooling equilibria.

<sup>6</sup>Reinforcement with iterative punishment and self-tuning reinforcement with iterative punishment are original to this paper.

<sup>7</sup>On the games and learning dynamics we will consider, the players will typically evolve *nearly* optimal dispositions on a successful run. Their dispositions will correspond to a bijective map from states to acts, but even when they are successful, their behavior will be somewhat noisy and hence not always respect the bijective map.


 FIGURE 1.  $4 \times 4 \times 4$  signaling game

Let  $q_i(t)$  be an agent's propensity for strategy  $i$  at time  $t$ . Her propensities evolve according to the update rule:

$$q_i(t+1) = \begin{cases} q_i(t) + \pi(t) & \text{if action } i \text{ was taken} \\ q_i(t) & \text{otherwise} \end{cases}$$

Here  $\pi(t)$  is the payoff received by an agent taking the action  $i$  on round  $t$ . The agent's current propensities determine her dispositions by determining the probability of a given action on each round. How this works is given by the response rule.

$$p_i(t) = \frac{q_i(t)}{\sum_j q_j(t)}$$

Here  $p_i(t)$  is the probability that the agent takes action  $i$  on round  $t$ .

Finally, one must also specify the initial propensities in order to say how the process gets started. Lower initial weights allow for more early exploration. Higher initial weights slow learning but make it more stable. And uneven initial weights bias early dispositions. As a consequence, different initial assignments may lead to very different behaviors. While one can assign these values any way one wants, we will suppose that each strategic option is given an equal and small initial weight, say  $q_i(0) = 1$  for all  $i$ . Here we will suppose that  $\pi(t) = 1$  on success and  $\pi(t) = 0$  otherwise.

On this dynamics, one might think of the players' dispositions and how they are updated in terms of urns with balls. On each play of the game, the current state of nature is determined in a random and unbiased way. The sender has one urn for each possible state of nature 0, 1, 2, 3. Each of these urns starts with one ball of each possible signal  $a, b, c, d$ . To determine the signal, she observes the current state of nature then draws a signal ball at random from the corresponding urn. She sends that signal. The receiver has one urn for each possible signal  $a, b, c, d$ . Each of these urns starts with one ball of each possible act 0, 1, 2, 3, where each state of nature requires the corresponding act for success. The receiver sees the signal then draws an act ball at random from the corresponding urn. She performs the act. If the receiver's act matches the current state, then each agent returns their ball to the urn from which it was drawn and adds a ball of the same type to the urn; otherwise, they simply replace return the ball that they drew.

The question as to how well a learning dynamics does on a particular type of problem is an empirical one. On simulation, when the modeled sender and receiver update their conditional dispositions under this dynamics, they start off randomly signaling and randomly acting. The signals initially have no meanings and there

is no pattern to the receiver’s action. But, as they learn from experience, about 0.797 of the time the players dispositions evolve to exhibit a cumulative success rate exceeding 0.8 with  $10^6$  plays per run (with 1000 runs). If they reach this level of success, then the run did not get stuck in a suboptimal pooling equilibrium as the most successful pooling equilibria for this game have cumulative success rates of 0.75. More generally, if a game evolves to do better than its best suboptimal pooling equilibria, then it tends to do monotonically better over time thereafter approaching an optimal signaling system. For this reason, the level of success that matters in each of the games we consider is the success rate of the most successful suboptimal equilibria. As for the specific case at hand, the players usually evolve a nearly optimal signaling system under simple reinforcement learning in the  $4 \times 4 \times 4$  game.

Basic reinforcement learning does not do as well on more complicated signaling games. With  $10^7$  plays per run, the  $8 \times 8 \times 8$  game exhibits a cumulative success rate over 0.875 just under half of the time 0.426. The higher cutoff here is because the most successful suboptimal pooling equilibria are better here. This is part of the reason the success rate is lower than for the  $4 \times 4 \times 4$  game. But the game is also harder. Note that we are running the  $8 \times 8 \times 8$  games ten times as long and still getting a significantly lower cumulative success rate.

While there is no entirely straightforward way of comparing the difficulty of evolving optimal dispositions for more and less complex games or of comparing the relative efficacy of different learning strategies, we will aim to give a sense of the relative difficulties and efficacies by describing what happens on simulation under an assortment of different game parameters. There are a few simplifying conventions that will allow for more accurate comparisons across different games. The first of these has to do with how we track success.

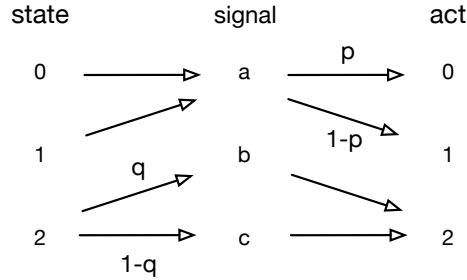
The most common way to measure evolutionary success in the signaling game literature is to track cumulative success over a fixed number of plays. The statistics reported above are for that.<sup>8</sup> An advantage of this measure is that one gets a rough sense of both how reliable the system ultimately is and how quickly it got there. In contrast, one might track the success rate over just the last few plays of the game on a run. This measure provides a better sense of the properties of the final evolved system. For the  $8 \times 8 \times 8$  game we just considered, 0.592 of the runs exhibit a success rate over 0.875 when just the last  $10^4$  plays of each run are considered. Given the focus of the present investigation, we will measure success in this way in what follows.

When the agents do not evolve toward an optimal signaling system, they get stuck in a suboptimal pooling equilibrium. As a concrete example of what this looks like, the suboptimal equilibrium illustrated in Figure 2 sometimes arises under simple reinforcement learning in the  $3 \times 3 \times 3$  game even with unbiased nature.

An optimal signaling system requires that one have a bijection between states of nature and signals and between signals and acts and hence, taken together, between states and signals. But here the sender uses  $a$  to represent both states 0 and 1 and the receiver randomly does act 0 and 1 with probabilities  $p$  and  $1 - p$  when he sees an  $a$  signal. And the sender randomly sends  $b$  and  $c$  with probabilities  $q$  and  $1 - q$  and the receiver always does act 2 when he gets either of these two signals.

---

<sup>8</sup>See Barrett (2006) for further examples of this measure for various signaling games on simple reinforcement learning and reinforcement with punishment.


 FIGURE 2.  $3 \times 3 \times 3$  signaling game

If nature is unbiased, the mean success rate will be  $2/3$ . And since unilaterally changing the value of either  $p$  or  $q$  does nothing to improve the agents' fortunes, their mixed strategies form an equilibrium under the dynamics.

A list of simulations results for various games under each of the forms of reinforcement learning that we discuss here can be found in section 6 at the end of the paper. We will discuss a number of these games in the text to give a sense of their behavior.

First, suboptimal equilibria are increasingly common the more complex the game. On  $10^6$  plays per run the  $16 \times 16 \times 16$  game has an end-of-run success rate on simple reinforcement learning better than 0.9375 on 0.169 of the runs. The  $32 \times 32 \times 32$  game does better than its most successful pooling equilibria at 0.969 on 0.003 of the runs. With ten times the number of plays per run  $10^7$ , the  $32 \times 32 \times 32$  game meets this cutoff a slightly more often on 0.036 of the runs. But on even  $2 \times 10^7$  plays per run,  $64 \times 64 \times 64$  game was never found to meet the corresponding cutoff at 0.985 on 1000 runs.

A natural question is whether one might do better on such complex games under a different learning dynamics. In one sense of better this question is easy to answer. There are some learning dynamics that will with probability one eventually evolve a strictly optimal signaling system for any finite signaling game. Consider win-stay/lose-randomize. Here the sender adopts an initial bijective map from states to signals and the receiver adopts an initial bijective map from signals to acts. On each play of the game, the sender and receiver stay with their present map unless the play fails in which case, each adopts a new bijective map independently and at random. Since there is a fixed positive probability on each play that they will happen to pick a pair of maps that fit with each other to produce a bijection between states and acts, this dynamics will eventually yield a optimal signaling system with probability one. But given the number of possible strategies in a moderately complex signaling game, it could take a cosmological long time to do so. In terms of speed this often a terrible dynamics. We will consider an illustrative example of this in the next section.

We turn now to consider how punishment might improve both the speed and accuracy of reinforcement learning in the context of signaling games.<sup>9</sup> This requires

<sup>9</sup>See Erev and Roth (1998), Barrett and Zollman (2009), Alexander, Skyrms, and Zabell (2012), Barrett, Cochran, Huttegger, and Fujiwara (2017), and Cochran and Barrett (2020) and (2021) for other promising learning strategies.

some care. While one can sometimes get a very fast and accurate learning dynamics by tuning the dynamics' parameters to the particular problem at hand, one should ultimately want a learning dynamics that does well in a broad range of contexts. We will aim to make progress on characterizing such a robust learning dynamics.

### 3. SIMPLE REINFORCEMENT WITH PUNISHMENT

A signaling game under simple reinforcement with punishment works very much like a signaling game under reinforcement except that one may also punish by removing balls from urns on failure. On  $(+i, -j)$  reinforcement with punishment, and agent's propensities evolve according to the update rule:

$$q_k(t+1) = \begin{cases} q_k(t) + i & \text{if action } k \text{ was taken and it was successful} \\ \max\{q_k(t) - j, 1\} & \text{if action } k \text{ was taken and it was unsuccessful} \\ q_k(t) & \text{otherwise} \end{cases}$$

The probability for each action is determined from the propensities in the usual way.

Consider an concrete example. In the  $4 \times 4 \times 4$  signaling game under  $(+i, -j)$  reinforcement with punishment, one adds  $i$  balls on success and may remove up to  $j$  balls on failure. More precisely, if the receiver's act matches the current state, then each agent returns their ball to the urn from which it was drawn and adds  $i$  balls of that type to the urn; otherwise, they return the ball to the urn from which it was drawn and remove  $j$  balls of that type from that urn subject to the condition that at least one ball of the type remains in the urn. Hence, if there are  $j$  or fewer balls of the given type, all are removed but one. The aim is to keep the same number of types in each urn by not allowing any type to go extinct.

Optimal signaling often evolves much faster and more reliably under reinforcement with punishment than under reinforcement alone. On just reinforcement learning without punishment, the  $8 \times 8 \times 8$  game exceeds the 0.875 final-success-rate cutoff on 0.613 of runs with  $10^6$  plays per run. Most of the failed runs here reflect suboptimal pooling equilibria. In contrast,  $(+1, -1)$  reinforcement with punishment exceeds the 0.875 success rate cutoff on a remarkable 0.945 of runs with the same number of plays per run.<sup>10</sup>

Punishment makes suboptimal pooling less likely by tending to weaken the players' suboptimal dispositions when they use them and fail. But it is also easy to punish too much. The  $8 \times 8 \times 8$  game with  $(+1, -2)$  learning was never found to exceed the 0.875 cutoff on 1000 runs with  $10^6$  plays per run. That said, increasing the level of reinforcement a bit preserves success for this particular game with 0.951 of runs meeting the cutoff for  $(+2 - 1)$  learning.

While  $(+1, -1)$  learning works well in the  $8 \times 8 \times 8$  game, it is less successful in more complex games. The  $16 \times 16 \times 16$  game was never observed to exceed the 0.9375 cutoff on  $(+1, -1)$  learning with  $10^6$  plays per run and 1000 runs. The coordination problem posed by the game is harder here. As a result, that level of punishment is too high relative to the level of reinforcement to allow for the effective evolution of optimal dispositions. Successful dispositions cannot get traction against the high level of punishment here. But the  $16 \times 16 \times 16$  game does exceed the cutoff 0.952

<sup>10</sup>While the form of reinforcement with punishment that we consider here has been studied elsewhere (see for example Barrett 2006), all of the simulation results from this section on are original to this paper.

of the time on  $(+2, -1)$  learning with  $10^6$  plays per run. Here the reinforcement on success is enough for the players to get traction on optimal dispositions early in the evolution. Indeed, the evolution of near optimal dispositions for these learning parameters is very fast. The cutoff is exceeded on 0.903 of the runs with just  $10^5$  plays per run.

While  $(+2, -1)$  learning works well in these two games it does not work at all in more complex games. The  $32 \times 32 \times 32$  game was never found to exceed a 0.969 cut off with this dynamics on 1000 runs with  $10^6$  plays per run. But this game did very well with  $(+3, -1)$  learning. Here it exceeded the cutoff on 0.879 of the runs. Again, one needs sufficient reinforcement relative to punishment to get traction early in the run. And success is highly sensitive to both parameter settings. Raising the level of punishment even slightly can make a striking difference. In the present game,  $(+3, -2)$  learning was never observed to succeed.

Similar phenomena are seen for yet more complex games. The  $64 \times 64 \times 64$  game was never observed to exceed a 0.984375 cutoff on  $(+3, -1)$  learning on 1000 runs with  $2 \times 10^6$  plays per run. But it was observed to exceed this strict cutoff on 0.604 of the runs with  $(+4, -1)$  learning. Indeed, it appears that the only thing limiting its success here was the relatively low number of plays per run given the complexity of the game as  $(+4, -1)$  learning did better than the just slightly lower 0.969 cutoff on 0.953 of the runs even with  $2 \times 10^6$  plays per run.

Determining what is happening in such simulations requires care. The  $128 \times 128 \times 128$  game was never observed to exceed a 0.995 cutoff on  $(+5, -1)$  learning with  $2 \times 10^6$  plays per run. But again it appears that the issue here was just one of run length given the high complexity of the game and high cutoff. This is strongly suggested by the fact that 0.780 of the runs exceed the cutoff with  $3 \times 10^7$  plays per run and 0.830 do with  $5 \times 10^7$  plays per run for the same reinforcement and punishment parameters.

The recurring theme here is that if the level of punishment is too low given the complexity of the game and the level of reinforcement, then one tends to get stuck in suboptimal pooling equilibria; and if it is too high given the complexity of the game and the level of reinforcement, then one cannot get the traction needed to evolve optimal dispositions. But both of these points also require care.<sup>11</sup>

On reinforcement learning with punishment, it is always logically possible to escape from a suboptimal pooling equilibrium. This might happen if one is lucky enough to fail, and hence weaken the dispositions that led to that failure, whenever one uses suboptimal dispositions. And if one is lucky enough, it is always logically possible to get the traction necessary to evolve optimal dispositions with even a very high level of punishment. A closer look at this second point will help to illustrate the role of tuning in simple reinforcement with punishment.

Consider a learning dynamics with a high level of punishment relative to the level of reinforcement. Regardless of how complex the game may be, given the randomly determined states of nature, it is in principle possible that they get lucky on the random dynamics that determines their actions and always succeed in the early plays of the game. If so, the punishment level does not matter early in the

<sup>11</sup>A similar point can be made with respect to reinforcement values. While  $(+5, -1)$  learning in the  $128 \times 128 \times 128$  game met its cutoff in 0.830 of runs (with  $5 \times 10^7$  plays per run),  $(+2, -0.4)$  never met the cutoff in the same game with the same number of plays per run. Just as when punishment is too high one cannot get the traction needed to evolve optimal dispositions, when reinforcement is too low one may not get the traction needed to evolve optimal dispositions.

run since the players will evolve nearly optimal dispositions by chance; and it does not matter later in the run either since the players will then have nearly optimal dispositions and hence rarely fail and be punished.

But while it is always *possible* for optimal signaling to evolve with even very high levels of punishment, one should expect simple reinforcement with punishment to be *extremely slow* with high levels of punishment relative to reinforcement, and particularly so for more complex games. In this, simple reinforcement with severe punishment behaves much like the win-stay/lose-randomize strategy discussed earlier.<sup>12</sup>

Given this, one might put the recurring theme in a somewhat more accurate positive form: simple reinforcement with punishment allows for the effective evolution of nearly optimal dispositions in even complex signaling games if one uses *just right* levels of reinforcement and punishment for the complexity of the game. It is worth reflecting on how remarkable this success may be before worrying further over tuning.

Setting aside the potentially unbounded number of mixed strategies available to the two players in a  $128 \times 128 \times 128$  game, there are  $128^{128}$  (more than  $10^{269}$ ) maps from states of nature to signals and the same number from signals to acts. But the vast majority of these maps are suboptimal. An optimal signaling system requires the sender to have a bijective map from states of nature to signals and the receiver to have a matching bijective map from signals to acts. There are  $128!$  (more than  $10^{215}$ ) bijective maps from states of nature to signals for the sender and the same number of bijective maps from signals to acts for the receiver. For a given sender bijective map, there is *only one* receiver map that will allow for optimal signaling. As a consequence, evolving optimal dispositions for such a game would be extraordinarily unlikely on something like win-stay/lose-randomize at supercomputing speeds and cosmological time scales even if the players strategies are restricted to just bijective maps.<sup>13</sup> That  $(+5 - 1)$  learning usually succeeds in finding a bijective map for the sender and a matching map for the receiver in just  $5 \times 10^7$  plays is a nontrivial virtue of simple reinforcement with punishment.

But while simple reinforcement with punishment may be remarkably effective with just-right learning parameters, one should want a *single* learning dynamics that does well on a broad assortment of games. To formulate something like this, we will start by considering reinforcement with iterative punishment. We will then reflect on how this dynamics might allow the players to tune how they learn as they learn to meet the demands of the signaling game at hand.

#### 4. REINFORCEMENT WITH ITERATIVE PUNISHMENT

On reinforcement with iterative punishment, one alternates between periods where the learning dynamics is mostly reinforcing and periods where it is mostly punishing. During a reinforcement phase, the players will form dispositions that together map from states of nature to receiver acts. Some parts of the map may be nearly optimal. On these a state of nature will produce the corresponding act with

<sup>12</sup>As we will see in the next section, reinforcement with iterative punishment allows for success with much higher levels of punishment since the early evolution is protected from punishment. Indeed, this will prove to be the key to the remarkable effectiveness of reinforcement with iterative punishment.

<sup>13</sup>The universe is likely younger than  $5 \times 10^{17}$  seconds.



high probability. Other parts may be suboptimal. On these the state may produce the corresponding action, but it does not do so reliably. A subsequent punishment phase will act to reset the dispositions associated with in the suboptimal mixed parts of the map, since their lower expected success rate is associated with more expected punishment, without undoing the nearly optimal bijective parts. Then the players will have another chance to get the suboptimal parts that were just erased right when they return to reinforcement. Inasmuch as they need only build the remaining parts of the bijective map on subsequent reinforcements, the players face an ever easier task after each punishment phase and hence stand a yet better chance of getting things right. When successful, they evolve nearly optimal dispositions by assembling the full bijective map in parts.

Tuning still matters for reinforcement with iterated punishment. One needs a sufficient level of reinforcement and needs to run it sufficiently long to get strong enough dispositions for the bijective parts of the map that it is unlikely that they will be erased by subsequent punishment. And one needs a sufficient level punishment and needs to run it long enough that it is likely that the parts of the map that are not bijective are erased. But one does not want the punishment to be so high or so long as to erase the nearly optimal bijective parts of the map that were hard-won on previous reinforcements.

Tuning in reinforcement with iterative punishment and in simple reinforcement with punishment are, however, quite different. In short, reinforcement with iterative punishment is much more robust. This is in part because it allows for a higher level of punishment. And this has to do with the core idea behind the dynamics.

Reinforcement with iterative punishment builds strong dispositions during the reinforcement phase, then uses punishment to erase the suboptimal parts of the map. Since the dispositions associated with the optimal parts of the bijective map are firmly established during reinforcement, even severe subsequent punishment will tend to leave these successful dispositions intact. So the players learn freely and without constraint during the reinforcement phase, and while they may evolve bad habits along the way, these are easily eliminated by severe subsequent punishment. In contrast, on simple reinforcement with punishment, punishment is always present and hence occurs before successful dispositions have been firmly established. As a result, weak dispositions that might ultimately form part of a successful bijective map if they were allowed to strengthen may be accidentally erased before they can prove their worth if the learning parameters are not just right. A high level of punishment makes it difficult for the players to build any dispositions at all. And a low level of punishment does not address the suboptimal pooling problem.

We will consider a few concrete examples of how reinforcement with iterative punishment works. While the reinforcement phases might in principle involve some degree of punishment and the punishment phases might involve some degree of reinforcement, we will consider here only dynamics with periods of pure reinforcement punctuated by periods of pure punishment. Using square brackets to represent reinforcement with iterative punishment,  $[+i, -j]$  learning alternates between periods of pure reinforcement of magnitude  $i$  and periods of pure punishment of magnitude  $j$ . Specifically, the agent's propensities evolve according to the update rule:

$$q_k(t+1) = \begin{cases} q_k(t) + i & \text{if action } k \text{ was taken and it was successful} \\ q_k(t) & \text{otherwise} \end{cases}$$

during the  $N$  plays of the reinforcement phase. Then they evolve according to the update rule:

$$q_k(t+1) = \begin{cases} \max\{q_k(t) - j, 1\} & \text{if action } k \text{ was taken and it was unsuccessful} \\ q_k(t) & \text{otherwise} \end{cases}$$

during the  $M$  plays of the punishment phase. Together the two phases constitute a block of play. The probability for each action is determined from the propensities in the usual way. The lengths of the periods  $N$  and  $M$  matter. We note these in each case.

Consider the  $32 \times 32 \times 32$  game under  $[+2, -4]$  reinforcement with iterative punishment. The first thing to note is that under  $(+2, -4)$  *simple* reinforcement with punishment, this level of punishment is too high for the players to get traction on a successful set of dispositions and they get nowhere. Indeed, as we saw earlier, even  $(+2, -1)$  simple reinforcement with punishment involves too much punishment for this complex a game. The players were never found to exceed the 0.969 success-rate cutoff. But  $[+2, -4]$  reinforcement with *iterative* punishment may do very well on this game depending on the length of the reinforcement and punishment phases. With a fixed run length of  $10^6$  plays per run and two iterations (one reinforces for a quarter of the run, then punishes for a quarter of the run, then repeats the reinforcement and punishment), the players exceed this high success-rate cutoff 0.252 of the time. With the same run length and four iterations, they exceed the cutoff 0.409 of the time. With eight iterations 0.645. And with sixteen iterations 0.918. The players are able to evolve nearly optimal dispositions very quickly on sixteen reinforcement/punishment blocks because the high punishment level erases any evolved suboptimal dispositions, and the high level of punishment does not undermine their optimal dispositions because it only kicks in after those dispositions are firmly established. This illustrates the radically different properties of reinforcement with iterative punishment. It also shows how the players build the full bijective map in parts on iterated punishment.

Importantly, the players do not continue to do better with more iterations with the fixed total run length. Alternating quickly between pure  $+2$  reinforcement and pure  $-4$  punishment closely approximates  $(+2, -4)$  simple reinforcement with punishment where the players are not at all effective in establishing successful dispositions.

In the same setup, the  $32 \times 32 \times 32$  game evolves more slowly on  $[+1, -2]$  than on  $[+2, -4]$  reinforcement with iterative punishment. With a fixed run length of  $10^6$  plays per run and two iterations, the players exceed the 0.969 success-rate cutoff 0.028 of the time. With the same run length and four iterations, they exceed the cutoff 0.050 of the time. With eight iterations 0.103. And with sixteen iterations 0.203. Part of the problem here is that the level of punishment is too low to fully erase the bad habits formed over a long period of reinforcement. Another is that we are not allowing the dynamics enough chances to get things right.

The  $32 \times 32 \times 32$  game on  $[+1, -2]$  reinforcement with iterative punishment provides a good example of the core idea behind reinforcement with iterative punishment and why it is robust. Consider this dynamics with a fixed reinforcement/punishment block length of  $1.25 \times 10^5$  plays. Eight such blocks exceeds the high 0.969 success-rate cutoff 0.008 of the time. Sixteen blocks exceeds the cutoff 0.323 of the time. Twenty-four blocks exceeds it 0.906 of the time. And thirty-two blocks 0.957 of the time. Here the players are observed to improve monotonically

with more reinforcement/punishment blocks. In each block sequence, the reinforcement phase builds potentially optimal dispositions and the punishment phase tends to erase suboptimal dispositions. The block length here is long enough to forge reasonably strong dispositions but short enough that the punishment phases accomplish significant cleanup. And repeating the process allows for incremental improvement.

Among the  $[+i, -j]$  parameters that were observed to do well on  $10^6$  plays per run are  $[+2, -3]$  (run success rate 0.839 with sixteen blocks),  $[+2, -6]$  (run success rate 0.837 with eight blocks), and  $[+2, -10]$  (run success rate 0.837 with four blocks). This illustrates the robustness of the dynamics under very different  $[+i, -j]$  parameters. It also illustrates how the dynamics with fixed  $[+i, -j]$  parameters may be tuned to the complexity of the game at hand by varying block length alone. There is further evidence of this latter sort of robustness.

Under reinforcement with iterative punishment, precisely the same  $[+i, -j]$  parameters, with a fixed block length, work well for games of very different complexity. On  $[+4, -12]$  learning and a block length of  $10^6$ , the  $32 \times 32 \times 32$  game meets its cutoff on 0.623 of the runs after  $3 \times 10^7$  plays per run; the  $64 \times 64 \times 64$  game meets its cutoff on 0.465 of the runs after  $5 \times 10^7$  plays per run;<sup>14</sup> the  $128 \times 128 \times 128$  game meets its cutoff on 0.903 of the runs after  $5 \times 10^7$  plays per run.<sup>15</sup> And in all of the runs for each of these three games (1000 runs each) the final success rate was observed to be better than 0.9. This last point indicates that the dynamics avoids both the suboptimal-pooling and the traction-under-high-punishment problems here. In doing so, it exhibits a radically different behavior from simple reinforcement with punishment, where the  $(+i, -j)$  parameters must be carefully tuned to the complexity of the signaling game.

So how might one choose block lengths that allow for monotonic improvement toward an optimal signaling system on iteration? Given the structure of reinforcement with iterated punishment, this can be done dynamically. What one wants is to run the reinforcement phase until one sees diminishing returns, then run the punishment phase until one sees diminishing returns, then repeat. Diminishing returns in the case of reinforcement means that one's local success rate is no longer improving significantly on repeated play. When this happens one typically has a some firmly established optimal dispositions with some mixed suboptimal dispositions. For punishment it means that one's local success rate is beginning to decline significantly. When this happens, one has erased the suboptimal dispositions and is beginning to erase optimal dispositions. Given the robustness of reinforcement with iterative punishment, one has considerable leeway in making these determinations.

Inasmuch as such determinations may be made during play, we have the framework for self-tuning reinforcement with iterative punishment. On this dynamics, one chooses a reinforcement level and a relatively severe punishment level, then one runs reinforcement until one sees diminishing returns, then runs punishment until one sees diminishing returns, then repeats. While there remains a significant

<sup>14</sup> $[+4, -12]$  performs even better with a block length that is tuned to the complexity of the game. On  $[+4, -12]$  learning with  $3 \times 10^7$  plays per run, the  $64 \times 64 \times 64$  game meets its cutoff on 0.694 and 0.912 of the runs (1000 runs each) with block lengths of  $5 \times 10^5$  and  $2.5 \times 10^5$  respectively.

<sup>15</sup>Cutoffs here are  $n - 1/n$ ; i.e. cutoff for the  $32 \times 32 \times 32$  game is 0.96875, for the  $64 \times 64 \times 64$  game it is 0.984375, and for the  $128 \times 128 \times 128$  game it is 0.9921875.

tuning issue in the choice of reinforcement and punishment levels, the dynamical determination of block length provides a significant degree of self-tuning.

There is more to say regarding how one might best detect diminishing returns to facilitate the rapid evolution of optimal dispositions. One should also want to consider how to choose reinforcement and punishment levels dynamically to this end. These are topics for further research.

## 5. DISCUSSION

While basic *reinforcement* learning allows for the evolution of optimal dispositions on simple signaling games, it often leads to suboptimal pooling equilibria, particularly in more complex games. *Simple reinforcement with punishment* helps to prevent suboptimal pooling, but it encounters a tuning problem. If the punishment level is too low, it does little to prevent suboptimal pooling; and if the punishment level is too high, successful dispositions fail to evolve at all.

In contrast, *reinforcement with iterated punishment* is tolerant of high levels of punishment. Here reinforcement establishes a mixed set of dispositions, and subsequent punishment tends to erase those dispositions that contribute to the suboptimal parts of the map from states to signals to acts. The aspect of signaling games that makes this dynamics effective is that an optimal bijective map can be built incrementally. The dynamics will move toward optimal dispositions if the reinforcement phase is long enough to accumulate firmly established dispositions and the punishment phase is long enough to erase the suboptimal dispositions but short enough not to do significant damage to the bijective parts of the map.

Understanding how reinforcement with iterative punishment works suggests *self-tuning reinforcement with iterated punishment*. Here one reinforces until progress begins to stall, punishes until the overall success of the system begins to decline, then repeats. While this still requires one to choose reinforcement and punishment levels up front, the players themselves determine the length of reinforcement/punishment blocks dynamically and hence have significant control of how much reinforcement and punishment in fact takes place. As a result, this dynamics is self-tuning over a broad class of signaling games.

Reinforcement with iterative punishment is more robust than simple reinforcement with punishment in two closely related ways. First, very different  $[+i, -j]$  parameters are found to evolve optimal dispositions in the same game. And second, precisely the same  $[+i, -j]$  parameters are found to evolve optimal dispositions in very different games. It is this that allows for the success of self-tuning reinforcement with iterative punishment. This self-tuning form of reinforcement learning clearly deserves further study.

Reinforcement with iterative punishment mirrors a learning structure that is often observed in natural contexts. During reinforcement periods, players learn quickly and establish a rich set of dispositions. Such low-risk free play is often punctuated by the higher-risk application of what the players have learned. On this dynamics, while they may make progress toward an optimal set of dispositions during free play, they may also pick up some bad habits along the way and end up with suboptimal dispositions. If higher-risk application tends to eliminate the suboptimal dispositions and if this happens in way that does not undermine the positive fruits of free play, subsequent free play allows for new, potentially optimal,

dispositions to be forged in the place of the bad habits that did not stand the test of real-world application.<sup>16</sup>

## 6. TABLES

In the tables below results are averaged over 1000 runs with the exception that the 128x128x128 iterated, 256x256x256 basic reinforcement with punishment, and 256x256x256 iterated games are averaged over 250 runs, and in the final table the 64x64x64 hybrid iterated game results are averaged over 200 runs. In the tables r1 and p1 are the reinforcement and punishment values respectively in the first half of an iteration while r2 and p2 are the reinforcement and punishment values for the second half of an iteration. When r1 = r2 and p1 = p2, then the learning dynamics are basic reinforcement with punishment; e.g. (+2,-1) is the same as r1 = +2, r2 = +2, p1 = -1, and p2 = -1. When r2 = p1 = 0, the dynamics are reinforcement with iterative punishment; e.g. [+2,-4] is the same as r1 = +2, r2 = +0, p1 = 0, and p2 = -4. The final two tables show results for a hybrid dynamics where r1 = r2, but p1  $\neq$  p2.

game	r 1	r 2	p 1	p 2	plays/run	run success rate
16x16x16	+1	+1	0	0	10 <sup>6</sup>	0.158
16x16x16	+1	+1	-1	-1	10 <sup>6</sup>	0.000
16x16x16	+1	+1	-2	-2	10 <sup>6</sup>	0.000
16x16x16	+1	+1	-3	-3	10 <sup>6</sup>	0.000
16x16x16	+1	+1	-4	-4	10 <sup>6</sup>	0.000
16x16x16	+1	+1	-5	-5	10 <sup>6</sup>	0.000
16x16x16	+2	+2	0	0	10 <sup>6</sup>	0.188
16x16x16	+2	+2	-1	-1	10 <sup>6</sup>	0.953
16x16x16	+2	+2	-2	-2	10 <sup>6</sup>	0.000
16x16x16	+2	+2	-3	-3	10 <sup>6</sup>	0.000
16x16x16	+2	+2	-4	-4	10 <sup>6</sup>	0.000
16x16x16	+2	+2	-5	-5	10 <sup>6</sup>	0.000
16x16x16	+3	+3	0	0	10 <sup>6</sup>	0.239
16x16x16	+3	+3	-1	-1	10 <sup>6</sup>	0.824
16x16x16	+3	+3	-2	-2	10 <sup>6</sup>	0.977
16x16x16	+3	+3	-3	-3	10 <sup>6</sup>	0.900
16x16x16	+3	+3	-4	-4	10 <sup>6</sup>	0.000
16x16x16	+3	+3	-5	-5	10 <sup>6</sup>	0.000
16x16x16	+4	+4	0	0	10 <sup>6</sup>	0.232
16x16x16	+4	+4	-1	-1	10 <sup>6</sup>	0.725
16x16x16	+4	+4	-2	-2	10 <sup>6</sup>	0.931
16x16x16	+4	+4	-3	-3	10 <sup>6</sup>	0.975
16x16x16	+4	+4	-4	-4	10 <sup>6</sup>	0.992
16x16x16	+4	+4	-5	-5	10 <sup>6</sup>	0.002
16x16x16	+5	+5	0	0	10 <sup>6</sup>	0.252
16x16x16	+5	+5	-1	-1	10 <sup>6</sup>	0.629
16x16x16	+5	+5	-2	-2	10 <sup>6</sup>	0.869
16x16x16	+5	+5	-3	-3	10 <sup>6</sup>	0.960
16x16x16	+5	+5	-4	-4	10 <sup>6</sup>	0.980

<sup>16</sup>We would like to thank Jack VanDrunen, Christian Torsell, and two anonymous reviewers for comments on an earlier version for this paper.

game	r 1	r 2	p 1	p 2	plays/run	run success rate
16x16x16	+5	+5	-5	-5	10 <sup>6</sup>	0.992
16x16x16	+6	+6	0	0	10 <sup>6</sup>	0.309
16x16x16	+6	+6	-1	-1	10 <sup>6</sup>	0.580
16x16x16	+6	+6	-2	-2	10 <sup>6</sup>	0.804
16x16x16	+6	+6	-3	-3	10 <sup>6</sup>	0.914
16x16x16	+6	+6	-4	-4	10 <sup>6</sup>	0.969
16x16x16	+6	+6	-5	-5	10 <sup>6</sup>	0.976
16x16x16	+7	+7	0	0	10 <sup>6</sup>	0.282
16x16x16	+7	+7	-1	-1	10 <sup>6</sup>	0.545
16x16x16	+7	+7	-2	-2	10 <sup>6</sup>	0.726
16x16x16	+7	+7	-3	-3	10 <sup>6</sup>	0.861
16x16x16	+7	+7	-4	-4	10 <sup>6</sup>	0.949
16x16x16	+7	+7	-5	-5	10 <sup>6</sup>	0.974
16x16x16	+8	+8	0	0	10 <sup>6</sup>	0.305
16x16x16	+8	+8	-1	-1	10 <sup>6</sup>	0.535
16x16x16	+8	+8	-2	-2	10 <sup>6</sup>	0.684
16x16x16	+8	+8	-3	-3	10 <sup>6</sup>	0.839
16x16x16	+8	+8	-4	-4	10 <sup>6</sup>	0.925
16x16x16	+8	+8	-5	-5	10 <sup>6</sup>	0.950
16x16x16	+9	+9	0	0	10 <sup>6</sup>	0.352
16x16x16	+9	+9	-1	-1	10 <sup>6</sup>	0.491
16x16x16	+9	+9	-2	-2	10 <sup>6</sup>	0.668
16x16x16	+9	+9	-3	-3	10 <sup>6</sup>	0.822
16x16x16	+9	+9	-4	-4	10 <sup>6</sup>	0.869
16x16x16	+9	+9	-5	-5	10 <sup>6</sup>	0.938
16x16x16	+10	+10	0	0	10 <sup>6</sup>	0.341
16x16x16	+10	+10	-1	-1	10 <sup>6</sup>	0.556
16x16x16	+10	+10	-2	-2	10 <sup>6</sup>	0.637
16x16x16	+10	+10	-3	-3	10 <sup>6</sup>	0.743
16x16x16	+10	+10	-4	-4	10 <sup>6</sup>	0.854
16x16x16	+10	+10	-5	-5	10 <sup>6</sup>	0.920

game	r 1	r 2	p 1	p 2	plays/run	run success rate
32x32x32	+1	+1	0	0	10 <sup>6</sup>	0.005
32x32x32	+1	+1	-1	-1	10 <sup>6</sup>	0.000
32x32x32	+1	+1	-2	-2	10 <sup>6</sup>	0.000
32x32x32	+1	+1	-3	-3	10 <sup>6</sup>	0.000
32x32x32	+1	+1	-4	-4	10 <sup>6</sup>	0.000
32x32x32	+1	+1	-5	-5	10 <sup>6</sup>	0.000
32x32x32	+2	+2	0	0	10 <sup>6</sup>	0.012
32x32x32	+2	+2	-1	-1	10 <sup>6</sup>	0.000
32x32x32	+2	+2	-2	-2	10 <sup>6</sup>	0.000
32x32x32	+2	+2	-3	-3	10 <sup>6</sup>	0.000
32x32x32	+2	+2	-4	-4	10 <sup>6</sup>	0.000
32x32x32	+2	+2	-5	-5	10 <sup>6</sup>	0.000
32x32x32	+3	+3	0	0	10 <sup>6</sup>	0.008
32x32x32	+3	+3	-1	-1	10 <sup>6</sup>	0.883

game	r 1	r 2	p 1	p 2	plays/run	run success rate
32x32x32	+3	+3	-2	-2	$10^6$	0.000
32x32x32	+3	+3	-3	-3	$10^6$	0.000
32x32x32	+3	+3	-4	-4	$10^6$	0.000
32x32x32	+3	+3	-5	-5	$10^6$	0.000
32x32x32	+4	+4	0	0	$10^6$	0.019
32x32x32	+4	+4	-1	-1	$10^6$	0.694
32x32x32	+4	+4	-2	-2	$10^6$	0.951
32x32x32	+4	+4	-3	-3	$10^6$	0.000
32x32x32	+4	+4	-4	-4	$10^6$	0.000
32x32x32	+4	+4	-5	-5	$10^6$	0.000
32x32x32	+5	+5	0	0	$10^6$	0.015
32x32x32	+5	+5	-1	-1	$10^6$	0.527
32x32x32	+5	+5	-2	-2	$10^6$	0.898
32x32x32	+5	+5	-3	-3	$10^6$	0.960
32x32x32	+5	+5	-4	-4	$10^6$	0.008
32x32x32	+5	+5	-5	-5	$10^6$	0.000
32x32x32	+6	+6	0	0	$10^6$	0.027
32x32x32	+6	+6	-1	-1	$10^6$	0.379
32x32x32	+6	+6	-2	-2	$10^6$	0.806
32x32x32	+6	+6	-3	-3	$10^6$	0.945
32x32x32	+6	+6	-4	-4	$10^6$	0.979
32x32x32	+6	+6	-5	-5	$10^6$	0.657
32x32x32	+7	+7	0	0	$10^6$	0.029
32x32x32	+7	+7	-1	-1	$10^6$	0.311
32x32x32	+7	+7	-2	-2	$10^6$	0.725
32x32x32	+7	+7	-3	-3	$10^6$	0.896
32x32x32	+7	+7	-4	-4	$10^6$	0.952
32x32x32	+7	+7	-5	-5	$10^6$	0.988
32x32x32	+8	+8	0	0	$10^6$	0.036
32x32x32	+8	+8	-1	-1	$10^6$	0.267
32x32x32	+8	+8	-2	-2	$10^6$	0.613
32x32x32	+8	+8	-3	-3	$10^6$	0.846
32x32x32	+8	+8	-4	-4	$10^6$	0.936
32x32x32	+8	+8	-5	-5	$10^6$	0.970
32x32x32	+9	+9	0	0	$10^6$	0.045
32x32x32	+9	+9	-1	-1	$10^6$	0.265
32x32x32	+9	+9	-2	-2	$10^6$	0.519
32x32x32	+9	+9	-3	-3	$10^6$	0.783
32x32x32	+9	+9	-4	-4	$10^6$	0.880
32x32x32	+9	+9	-5	-5	$10^6$	0.945
32x32x32	+10	+10	0	0	$10^6$	0.042
32x32x32	+10	+10	-1	-1	$10^6$	0.229
32x32x32	+10	+10	-2	-2	$10^6$	0.483
32x32x32	+10	+10	-3	-3	$10^6$	0.712
32x32x32	+10	+10	-4	-4	$10^6$	0.859
32x32x32	+10	+10	-5	-5	$10^6$	0.925

game	r 1	r 2	p 1	p 2	plays/run	run success rate
64x64x64	+1	+1	0	0	$10^7$	0.000
64x64x64	+1	+1	-1	-1	$10^7$	0.000
64x64x64	+1	+1	-2	-2	$10^7$	0.000
64x64x64	+1	+1	-3	-3	$10^7$	0.000
64x64x64	+1	+1	-4	-4	$10^7$	0.000
64x64x64	+1	+1	-5	-5	$10^7$	0.000
64x64x64	+2	+2	0	0	$10^7$	0.001
64x64x64	+2	+2	-1	-1	$10^7$	0.000
64x64x64	+2	+2	-2	-2	$10^7$	0.000
64x64x64	+2	+2	-3	-3	$10^7$	0.000
64x64x64	+2	+2	-4	-4	$10^7$	0.000
64x64x64	+2	+2	-5	-5	$10^7$	0.000
64x64x64	+3	+3	0	0	$10^7$	0.002
64x64x64	+3	+3	-1	-1	$10^7$	0.000
64x64x64	+3	+3	-2	-2	$10^7$	0.000
64x64x64	+3	+3	-3	-3	$10^7$	0.000
64x64x64	+3	+3	-4	-4	$10^7$	0.000
64x64x64	+3	+3	-5	-5	$10^7$	0.000
64x64x64	+4	+4	0	0	$10^7$	0.000
64x64x64	+4	+4	-1	-1	$10^7$	0.871
64x64x64	+4	+4	-2	-2	$10^7$	0.000
64x64x64	+4	+4	-3	-3	$10^7$	0.000
64x64x64	+4	+4	-4	-4	$10^7$	0.000
64x64x64	+4	+4	-5	-5	$10^7$	0.000
64x64x64	+5	+5	0	0	$10^7$	0.000
64x64x64	+5	+5	-1	-1	$10^7$	0.714
64x64x64	+5	+5	-2	-2	$10^7$	0.967
64x64x64	+5	+5	-3	-3	$10^7$	0.000
64x64x64	+5	+5	-4	-4	$10^7$	0.000
64x64x64	+5	+5	-5	-5	$10^7$	0.000
64x64x64	+6	+6	0	0	$10^7$	0.003
64x64x64	+6	+6	-1	-1	$10^7$	0.462
64x64x64	+6	+6	-2	-2	$10^7$	0.914
64x64x64	+6	+6	-3	-3	$10^7$	0.915
64x64x64	+6	+6	-4	-4	$10^7$	0.000
64x64x64	+6	+6	-5	-5	$10^7$	0.000
64x64x64	+7	+7	0	0	$10^7$	0.001
64x64x64	+7	+7	-1	-1	$10^7$	0.279
64x64x64	+7	+7	-2	-2	$10^7$	0.855
64x64x64	+7	+7	-3	-3	$10^7$	0.953
64x64x64	+7	+7	-4	-4	$10^7$	0.939
64x64x64	+7	+7	-5	-5	$10^7$	0.000
64x64x64	+8	+8	0	0	$10^7$	0.001
64x64x64	+8	+8	-1	-1	$10^7$	0.169
64x64x64	+8	+8	-2	-2	$10^7$	0.744
64x64x64	+8	+8	-3	-3	$10^7$	0.938
64x64x64	+8	+8	-4	-4	$10^7$	0.965
64x64x64	+8	+8	-5	-5	$10^7$	0.987
64x64x64	+9	+9	0	0	$10^7$	0.001



game	r 1	r 2	p 1	p 2	plays/run	run success rate
64x64x64	+9	+9	-1	-1	$10^7$	0.141
64x64x64	+9	+9	-2	-2	$10^7$	0.633
64x64x64	+9	+9	-3	-3	$10^7$	0.890
64x64x64	+9	+9	-4	-4	$10^7$	0.966
64x64x64	+9	+9	-5	-5	$10^7$	0.977
64x64x64	+10	+10	0	0	$10^7$	0.001
64x64x64	+10	+10	-1	-1	$10^7$	0.098
64x64x64	+10	+10	-2	-2	$10^7$	0.483
64x64x64	+10	+10	-3	-3	$10^7$	0.804
64x64x64	+10	+10	-4	-4	$10^7$	0.929
64x64x64	+10	+10	-5	-5	$10^7$	0.964

game	r 1	r 2	p 1	p 2	plays/run	run success rate
128x128x128	+1	+1	0	0	$2 \cdot 10^7$	0.000
128x128x128	+1	+1	-1	-1	$2 \cdot 10^7$	0.000
128x128x128	+1	+1	-2	-2	$2 \cdot 10^7$	0.000
128x128x128	+1	+1	-3	-3	$2 \cdot 10^7$	0.000
128x128x128	+1	+1	-4	-4	$2 \cdot 10^7$	0.000
128x128x128	+1	+1	-5	-5	$2 \cdot 10^7$	0.000
128x128x128	+2	+2	0	0	$2 \cdot 10^7$	0.000
128x128x128	+2	+2	-1	-1	$2 \cdot 10^7$	0.000
128x128x128	+2	+2	-2	-2	$2 \cdot 10^7$	0.000
128x128x128	+2	+2	-3	-3	$2 \cdot 10^7$	0.000
128x128x128	+2	+2	-4	-4	$2 \cdot 10^7$	0.000
128x128x128	+2	+2	-5	-5	$2 \cdot 10^7$	0.000
128x128x128	+3	+3	0	0	$2 \cdot 10^7$	0.000
128x128x128	+3	+3	-1	-1	$2 \cdot 10^7$	0.000
128x128x128	+3	+3	-2	-2	$2 \cdot 10^7$	0.000
128x128x128	+3	+3	-3	-3	$2 \cdot 10^7$	0.000
128x128x128	+3	+3	-4	-4	$2 \cdot 10^7$	0.000
128x128x128	+3	+3	-5	-5	$2 \cdot 10^7$	0.000
128x128x128	+4	+4	0	0	$2 \cdot 10^7$	0.000
128x128x128	+4	+4	-1	-1	$2 \cdot 10^7$	0.000
128x128x128	+4	+4	-2	-2	$2 \cdot 10^7$	0.000
128x128x128	+4	+4	-3	-3	$2 \cdot 10^7$	0.000
128x128x128	+4	+4	-4	-4	$2 \cdot 10^7$	0.000
128x128x128	+4	+4	-5	-5	$2 \cdot 10^7$	0.000
128x128x128	+5	+5	0	0	$2 \cdot 10^7$	0.000
128x128x128	+5	+5	-1	-1	$2 \cdot 10^7$	0.010
128x128x128	+5	+5	-2	-2	$2 \cdot 10^7$	0.000
128x128x128	+5	+5	-3	-3	$2 \cdot 10^7$	0.000
128x128x128	+5	+5	-4	-4	$2 \cdot 10^7$	0.000
128x128x128	+5	+5	-5	-5	$2 \cdot 10^7$	0.000
128x128x128	+6	+6	0	0	$2 \cdot 10^7$	0.000
128x128x128	+6	+6	-1	-1	$2 \cdot 10^7$	0.740
128x128x128	+6	+6	-2	-2	$2 \cdot 10^7$	0.000
128x128x128	+6	+6	-3	-3	$2 \cdot 10^7$	0.000

game	r 1	r 2	p 1	p 2	plays/run	run success rate
128x128x128	+6	+6	-4	-4	$2 \cdot 10^7$	0.000
128x128x128	+6	+6	-5	-5	$2 \cdot 10^7$	0.000
128x128x128	+7	+7	0	0	$2 \cdot 10^7$	0.000
128x128x128	+7	+7	-1	-1	$2 \cdot 10^7$	0.481
128x128x128	+7	+7	-2	-2	$2 \cdot 10^7$	0.000
128x128x128	+7	+7	-3	-3	$2 \cdot 10^7$	0.000
128x128x128	+7	+7	-4	-4	$2 \cdot 10^7$	0.000
128x128x128	+7	+7	-5	-5	$2 \cdot 10^7$	0.000
128x128x128	+8	+8	0	0	$2 \cdot 10^7$	0.000
128x128x128	+8	+8	-1	-1	$2 \cdot 10^7$	0.312
128x128x128	+8	+8	-2	-2	$2 \cdot 10^7$	0.868
128x128x128	+8	+8	-3	-3	$2 \cdot 10^7$	0.000
128x128x128	+8	+8	-4	-4	$2 \cdot 10^7$	0.000
128x128x128	+8	+8	-5	-5	$2 \cdot 10^7$	0.000
128x128x128	+9	+9	0	0	$2 \cdot 10^7$	0.000
128x128x128	+9	+9	-1	-1	$2 \cdot 10^7$	0.141
128x128x128	+9	+9	-2	-2	$2 \cdot 10^7$	0.801
128x128x128	+9	+9	-3	-3	$2 \cdot 10^7$	0.268
128x128x128	+9	+9	-4	-4	$2 \cdot 10^7$	0.000
128x128x128	+9	+9	-5	-5	$2 \cdot 10^7$	0.000
128x128x128	+10	+10	0	0	$2 \cdot 10^7$	0.000
128x128x128	+10	+10	-1	-1	$2 \cdot 10^7$	0.077
128x128x128	+10	+10	-2	-2	$2 \cdot 10^7$	0.691
128x128x128	+10	+10	-3	-3	$2 \cdot 10^7$	0.927
128x128x128	+10	+10	-4	-4	$2 \cdot 10^7$	0.000
128x128x128	+10	+10	-5	-5	$2 \cdot 10^7$	0.000

game	r 1	r 2	p 1	p 2	plays/run	run success rate
256x256x256	+11	+11	0	0	$5 \cdot 10^7$	0.000
256x256x256	+11	+11	-1	-1	$5 \cdot 10^7$	0.084
256x256x256	+11	+11	-2	-2	$5 \cdot 10^7$	0.000
256x256x256	+11	+11	-3	-3	$5 \cdot 10^7$	0.000
256x256x256	+11	+11	-4	-4	$5 \cdot 10^7$	0.000
256x256x256	+11	+11	-5	-5	$5 \cdot 10^7$	0.000
256x256x256	+12	+12	0	0	$5 \cdot 10^7$	0.000
256x256x256	+12	+12	-1	-1	$5 \cdot 10^7$	0.020
256x256x256	+12	+12	-2	-2	$5 \cdot 10^7$	0.532
256x256x256	+12	+12	-3	-3	$5 \cdot 10^7$	0.000
256x256x256	+12	+12	-4	-4	$5 \cdot 10^7$	0.000
256x256x256	+12	+12	-5	-5	$5 \cdot 10^7$	0.000
256x256x256	+13	+13	0	0	$5 \cdot 10^7$	0.000
256x256x256	+13	+13	-1	-1	$5 \cdot 10^7$	0.020
256x256x256	+13	+13	-2	-2	$5 \cdot 10^7$	0.572
256x256x256	+13	+13	-3	-3	$5 \cdot 10^7$	0.000
256x256x256	+13	+13	-4	-4	$5 \cdot 10^7$	0.000
256x256x256	+13	+13	-5	-5	$5 \cdot 10^7$	0.000
256x256x256	+14	+14	0	0	$5 \cdot 10^7$	0.000

game	r 1	r 2	p 1	p 2	plays/run	run success rate
256x256x256	+14	+14	-1	-1	$5 \cdot 10^{-7}$	0.004
256x256x256	+14	+14	-2	-2	$5 \cdot 10^{-7}$	0.476
256x256x256	+14	+14	-3	-3	$5 \cdot 10^{-7}$	0.232
256x256x256	+14	+14	-4	-4	$5 \cdot 10^{-7}$	0.000
256x256x256	+14	+14	-5	-5	$5 \cdot 10^{-7}$	0.000
256x256x256	+15	+15	0	0	$5 \cdot 10^{-7}$	0.000
256x256x256	+15	+15	-1	-1	$5 \cdot 10^{-7}$	0.008
256x256x256	+15	+15	-2	-2	$5 \cdot 10^{-7}$	0.368
256x256x256	+15	+15	-3	-3	$5 \cdot 10^{-7}$	0.488
256x256x256	+15	+15	-4	-4	$5 \cdot 10^{-7}$	0.000
256x256x256	+15	+15	-5	-5	$5 \cdot 10^{-7}$	0.000
256x256x256	+16	+16	0	0	$5 \cdot 10^{-7}$	0.000
256x256x256	+16	+16	-1	-1	$5 \cdot 10^{-7}$	0.004
256x256x256	+16	+16	-2	-2	$5 \cdot 10^{-7}$	0.280
256x256x256	+16	+16	-3	-3	$5 \cdot 10^{-7}$	0.624
256x256x256	+16	+16	-4	-4	$5 \cdot 10^{-7}$	0.048
256x256x256	+16	+16	-5	-5	$5 \cdot 10^{-7}$	0.000
256x256x256	+17	+17	0	0	$5 \cdot 10^{-7}$	0.000
256x256x256	+17	+17	-1	-1	$5 \cdot 10^{-7}$	0.000
256x256x256	+17	+17	-2	-2	$5 \cdot 10^{-7}$	0.184
256x256x256	+17	+17	-3	-3	$5 \cdot 10^{-7}$	0.604
256x256x256	+17	+17	-4	-4	$5 \cdot 10^{-7}$	0.568
256x256x256	+17	+17	-5	-5	$5 \cdot 10^{-7}$	0.000
256x256x256	+18	+18	0	0	$5 \cdot 10^{-7}$	0.000
256x256x256	+18	+18	-1	-1	$5 \cdot 10^{-7}$	0.000
256x256x256	+18	+18	-2	-2	$5 \cdot 10^{-7}$	0.136
256x256x256	+18	+18	-3	-3	$5 \cdot 10^{-7}$	0.488
256x256x256	+18	+18	-4	-4	$5 \cdot 10^{-7}$	0.704
256x256x256	+18	+18	-5	-5	$5 \cdot 10^{-7}$	0.092
256x256x256	+19	+19	0	0	$5 \cdot 10^{-7}$	0.000
256x256x256	+19	+19	-1	-1	$5 \cdot 10^{-7}$	0.000
256x256x256	+19	+19	-2	-2	$5 \cdot 10^{-7}$	0.076
256x256x256	+19	+19	-3	-3	$5 \cdot 10^{-7}$	0.436
256x256x256	+19	+19	-4	-4	$5 \cdot 10^{-7}$	0.708
256x256x256	+19	+19	-5	-5	$5 \cdot 10^{-7}$	0.588
256x256x256	+20	+20	0	0	$5 \cdot 10^{-7}$	0.000
256x256x256	+20	+20	-1	-1	$5 \cdot 10^{-7}$	0.000
256x256x256	+20	+20	-2	-2	$5 \cdot 10^{-7}$	0.076
256x256x256	+20	+20	-3	-3	$5 \cdot 10^{-7}$	0.384
256x256x256	+20	+20	-4	-4	$5 \cdot 10^{-7}$	0.648
256x256x256	+20	+20	-5	-5	$5 \cdot 10^{-7}$	0.632

In the following tables, “rsr  $x$  iters” is the run success rate when the length of the full run is broken into  $x$  iterations. The iteration length on a game will hence be different for different values of  $x$ .

game	r 1	r 2	p 1	p 2	plays/run	rsr 2 iters	rsr 4 iters	rsr 8 iters	rsr 16 iters
32x32x32	+1	+0	0	-1	$10^6$	0.136	0.306	0.684	0.684

game	r 1	r 2	p 1	p 2	plays/run	rsr 2 iters	rsr 4 iters	rsr 8 iters	rsr 16 iters
32x32x32	+1	+0	0	-2	10 <sup>6</sup>	0.253	0.644	0.008	0.008
32x32x32	+1	+0	0	-3	10 <sup>6</sup>	0.359	0.667	0.000	0.000
32x32x32	+1	+0	0	-4	10 <sup>6</sup>	0.293	0.078	0.000	0.000
32x32x32	+1	+0	0	-5	10 <sup>6</sup>	0.169	0.007	0.000	0.000
32x32x32	+1	+0	0	-6	10 <sup>6</sup>	0.064	0.000	0.000	0.000
32x32x32	+1	+0	0	-7	10 <sup>6</sup>	0.046	0.000	0.000	0.000
32x32x32	+1	+0	0	-8	10 <sup>6</sup>	0.019	0.000	0.000	0.000
32x32x32	+1	+0	0	-9	10 <sup>6</sup>	0.010	0.000	0.000	0.000
32x32x32	+1	+0	0	-10	10 <sup>6</sup>	0.010	0.000	0.000	0.000
32x32x32	+2	+0	0	-1	10 <sup>6</sup>	0.028	0.050	0.103	0.103
32x32x32	+2	+0	0	-2	10 <sup>6</sup>	0.104	0.138	0.252	0.252
32x32x32	+2	+0	0	-3	10 <sup>6</sup>	0.171	0.268	0.474	0.474
32x32x32	+2	+0	0	-4	10 <sup>6</sup>	0.252	0.409	0.645	0.645
32x32x32	+2	+0	0	-5	10 <sup>6</sup>	0.310	0.549	0.762	0.762
32x32x32	+2	+0	0	-6	10 <sup>6</sup>	0.360	0.661	0.837	0.837
32x32x32	+2	+0	0	-7	10 <sup>6</sup>	0.416	0.780	0.895	0.895
32x32x32	+2	+0	0	-8	10 <sup>6</sup>	0.479	0.808	0.882	0.882
32x32x32	+2	+0	0	-9	10 <sup>6</sup>	0.508	0.823	0.609	0.609
32x32x32	+2	+0	0	-10	10 <sup>6</sup>	0.447	0.830	0.273	0.273
32x32x32	+3	+0	0	-1	10 <sup>6</sup>	0.018	0.030	0.055	0.055
32x32x32	+3	+0	0	-2	10 <sup>6</sup>	0.044	0.094	0.161	0.161
32x32x32	+3	+0	0	-3	10 <sup>6</sup>	0.087	0.186	0.219	0.219
32x32x32	+3	+0	0	-4	10 <sup>6</sup>	0.142	0.211	0.280	0.280
32x32x32	+3	+0	0	-5	10 <sup>6</sup>	0.196	0.308	0.402	0.402
32x32x32	+3	+0	0	-6	10 <sup>6</sup>	0.236	0.334	0.527	0.527
32x32x32	+3	+0	0	-7	10 <sup>6</sup>	0.304	0.463	0.633	0.633
32x32x32	+3	+0	0	-8	10 <sup>6</sup>	0.306	0.504	0.734	0.734
32x32x32	+3	+0	0	-9	10 <sup>6</sup>	0.383	0.617	0.793	0.793
32x32x32	+3	+0	0	-10	10 <sup>6</sup>	0.413	0.698	0.847	0.847
32x32x32	+4	+0	0	-1	10 <sup>6</sup>	0.014	0.034	0.042	0.042
32x32x32	+4	+0	0	-2	10 <sup>6</sup>	0.027	0.056	0.101	0.101
32x32x32	+4	+0	0	-3	10 <sup>6</sup>	0.050	0.110	0.155	0.155
32x32x32	+4	+0	0	-4	10 <sup>6</sup>	0.095	0.184	0.235	0.235
32x32x32	+4	+0	0	-5	10 <sup>6</sup>	0.122	0.210	0.290	0.290
32x32x32	+4	+0	0	-6	10 <sup>6</sup>	0.176	0.291	0.380	0.380
32x32x32	+4	+0	0	-7	10 <sup>6</sup>	0.237	0.314	0.376	0.376
32x32x32	+4	+0	0	-8	10 <sup>6</sup>	0.238	0.398	0.496	0.496
32x32x32	+4	+0	0	-9	10 <sup>6</sup>	0.317	0.457	0.557	0.557
32x32x32	+4	+0	0	-10	10 <sup>6</sup>	0.349	0.516	0.636	0.636

game	r 1	r 2	p 1	p 2	plays/run	rsr 10 iters	rsr 20 iters	rsr 80 iters
64x64x64	+1	+0	0	-0	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+1	+0	0	-1	10 <sup>7</sup>	0.302	0.836	0.000
64x64x64	+1	+0	0	-2	10 <sup>7</sup>	0.675	0.000	0.000
64x64x64	+1	+0	0	-3	10 <sup>7</sup>	0.016	0.000	0.000
64x64x64	+1	+0	0	-4	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+1	+0	0	-5	10 <sup>7</sup>	0.000	0.000	0.000

game	r 1	r 2	p 1	p 2	plays/run	rsr 10 iters	rsr 20 iters	rsr 80 iters
64x64x64	+1	+0	0	-6	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+1	+0	0	-7	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+1	+0	0	-8	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+1	+0	0	-9	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+1	+0	0	-10	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+1	+0	0	-11	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+1	+0	0	-12	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+1	+0	0	-13	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+1	+0	0	-14	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+1	+0	0	-15	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+2	+0	0	-0	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+2	+0	0	-1	10 <sup>7</sup>	0.017	0.057	0.910
64x64x64	+2	+0	0	-2	10 <sup>7</sup>	0.074	0.247	0.000
64x64x64	+2	+0	0	-3	10 <sup>7</sup>	0.172	0.493	0.000
64x64x64	+2	+0	0	-4	10 <sup>7</sup>	0.281	0.674	0.000
64x64x64	+2	+0	0	-5	10 <sup>7</sup>	0.452	0.757	0.000
64x64x64	+2	+0	0	-6	10 <sup>7</sup>	0.582	0.859	0.000
64x64x64	+2	+0	0	-7	10 <sup>7</sup>	0.712	0.917	0.000
64x64x64	+2	+0	0	-8	10 <sup>7</sup>	0.742	0.227	0.000
64x64x64	+2	+0	0	-9	10 <sup>7</sup>	0.780	0.001	0.000
64x64x64	+2	+0	0	-10	10 <sup>7</sup>	0.807	0.000	0.000
64x64x64	+2	+0	0	-11	10 <sup>7</sup>	0.821	0.000	0.000
64x64x64	+2	+0	0	-12	10 <sup>7</sup>	0.755	0.000	0.000
64x64x64	+2	+0	0	-13	10 <sup>7</sup>	0.542	0.000	0.000
64x64x64	+2	+0	0	-14	10 <sup>7</sup>	0.258	0.000	0.000
64x64x64	+2	+0	0	-15	10 <sup>7</sup>	0.143	0.000	0.000
64x64x64	+3	+0	0	-0	10 <sup>7</sup>	0.000	0.000	0.001
64x64x64	+3	+0	0	-1	10 <sup>7</sup>	0.006	0.013	0.313
64x64x64	+3	+0	0	-2	10 <sup>7</sup>	0.026	0.051	0.771
64x64x64	+3	+0	0	-3	10 <sup>7</sup>	0.063	0.126	0.913
64x64x64	+3	+0	0	-4	10 <sup>7</sup>	0.100	0.214	0.962
64x64x64	+3	+0	0	-5	10 <sup>7</sup>	0.160	0.323	0.957
64x64x64	+3	+0	0	-6	10 <sup>7</sup>	0.223	0.439	0.000
64x64x64	+3	+0	0	-7	10 <sup>7</sup>	0.290	0.562	0.000
64x64x64	+3	+0	0	-8	10 <sup>7</sup>	0.379	0.678	0.000
64x64x64	+3	+0	0	-9	10 <sup>7</sup>	0.474	0.745	0.000
64x64x64	+3	+0	0	-10	10 <sup>7</sup>	0.562	0.800	0.000
64x64x64	+3	+0	0	-11	10 <sup>7</sup>	0.586	0.802	0.000
64x64x64	+3	+0	0	-12	10 <sup>7</sup>	0.694	0.841	0.000
64x64x64	+3	+0	0	-13	10 <sup>7</sup>	0.733	0.860	0.000
64x64x64	+3	+0	0	-14	10 <sup>7</sup>	0.799	0.900	0.000
64x64x64	+3	+0	0	-15	10 <sup>7</sup>	0.794	0.908	0.000
64x64x64	+4	+0	0	-0	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+4	+0	0	-1	10 <sup>7</sup>	0.002	0.009	0.050
64x64x64	+4	+0	0	-2	10 <sup>7</sup>	0.017	0.034	0.252
64x64x64	+4	+0	0	-3	10 <sup>7</sup>	0.033	0.062	0.521
64x64x64	+4	+0	0	-4	10 <sup>7</sup>	0.075	0.079	0.716
64x64x64	+4	+0	0	-5	10 <sup>7</sup>	0.097	0.127	0.840
64x64x64	+4	+0	0	-6	10 <sup>7</sup>	0.129	0.223	0.915

game	r 1	r 2	p 1	p 2	plays/run	rsr 10 iters	rsr 20 iters	rsr 80 iters
64x64x64	+4	+0	0	-7	10 <sup>7</sup>	0.149	0.258	0.947
64x64x64	+4	+0	0	-8	10 <sup>7</sup>	0.207	0.324	0.961
64x64x64	+4	+0	0	-9	10 <sup>7</sup>	0.230	0.433	0.978
64x64x64	+4	+0	0	-10	10 <sup>7</sup>	0.300	0.494	0.486
64x64x64	+4	+0	0	-11	10 <sup>7</sup>	0.346	0.578	0.000
64x64x64	+4	+0	0	-12	10 <sup>7</sup>	0.413	0.642	0.000
64x64x64	+4	+0	0	-13	10 <sup>7</sup>	0.481	0.684	0.000
64x64x64	+4	+0	0	-14	10 <sup>7</sup>	0.534	0.745	0.000
64x64x64	+4	+0	0	-15	10 <sup>7</sup>	0.589	0.791	0.000
64x64x64	+5	+0	0	-0	10 <sup>7</sup>	0.000	0.000	0.000
64x64x64	+5	+0	0	-1	10 <sup>7</sup>	0.005	0.004	0.010
64x64x64	+5	+0	0	-2	10 <sup>7</sup>	0.015	0.020	0.082
64x64x64	+5	+0	0	-3	10 <sup>7</sup>	0.028	0.032	0.177
64x64x64	+5	+0	0	-4	10 <sup>7</sup>	0.041	0.074	0.331
64x64x64	+5	+0	0	-5	10 <sup>7</sup>	0.069	0.091	0.527
64x64x64	+5	+0	0	-6	10 <sup>7</sup>	0.109	0.130	0.661
64x64x64	+5	+0	0	-7	10 <sup>7</sup>	0.135	0.167	0.770
64x64x64	+5	+0	0	-8	10 <sup>7</sup>	0.146	0.179	0.835
64x64x64	+5	+0	0	-9	10 <sup>7</sup>	0.186	0.236	0.879
64x64x64	+5	+0	0	-10	10 <sup>7</sup>	0.212	0.286	0.915
64x64x64	+5	+0	0	-11	10 <sup>7</sup>	0.245	0.353	0.934
64x64x64	+5	+0	0	-12	10 <sup>7</sup>	0.259	0.406	0.947
64x64x64	+5	+0	0	-13	10 <sup>7</sup>	0.281	0.471	0.964
64x64x64	+5	+0	0	-14	10 <sup>7</sup>	0.345	0.500	0.978
64x64x64	+5	+0	0	-15	10 <sup>7</sup>	0.407	0.596	0.980

game	r 1	r 2	p 1	p 2	plays/run	rsr 25 iters	rsr 50 iters
128x128x128	+2	+0	0	-5	2.5*10 <sup>7</sup>	0.000	0.000
128x128x128	+2	+0	0	-6	2.5*10 <sup>7</sup>	0.000	0.000
128x128x128	+2	+0	0	-7	2.5*10 <sup>7</sup>	0.000	0.000
128x128x128	+2	+0	0	-8	2.5*10 <sup>7</sup>	0.000	0.000
128x128x128	+2	+0	0	-9	2.5*10 <sup>7</sup>	0.000	0.000
128x128x128	+2	+0	0	-10	2.5*10 <sup>7</sup>	0.000	0.000
128x128x128	+2	+0	0	-11	2.5*10 <sup>7</sup>	0.000	0.000
128x128x128	+2	+0	0	-12	2.5*10 <sup>7</sup>	0.000	0.000
128x128x128	+2	+0	0	-13	2.5*10 <sup>7</sup>	0.000	0.000
128x128x128	+2	+0	0	-14	2.5*10 <sup>7</sup>	0.000	0.000
128x128x128	+2	+0	0	-15	2.5*10 <sup>7</sup>	0.000	0.000
128x128x128	+4	+0	0	-5	2.5*10 <sup>7</sup>	0.204	0.920
128x128x128	+4	+0	0	-6	2.5*10 <sup>7</sup>	0.464	0.848
128x128x128	+4	+0	0	-7	2.5*10 <sup>7</sup>	0.512	0.000
128x128x128	+4	+0	0	-8	2.5*10 <sup>7</sup>	0.736	0.000
128x128x128	+4	+0	0	-9	2.5*10 <sup>7</sup>	0.732	0.000
128x128x128	+4	+0	0	-10	2.5*10 <sup>7</sup>	0.804	0.000
128x128x128	+4	+0	0	-11	2.5*10 <sup>7</sup>	0.852	0.000
128x128x128	+4	+0	0	-12	2.5*10 <sup>7</sup>	0.892	0.000
128x128x128	+4	+0	0	-13	2.5*10 <sup>7</sup>	0.880	0.000

game	r 1	r 2	p 1	p 2	plays/run	rsr 25 iters	rsr 50 iters
128x128x128	+4	+0	0	-14	$2.5 \cdot 10^{-7}$	0.444	0.000
128x128x128	+4	+0	0	-15	$2.5 \cdot 10^{-7}$	0.008	0.000
128x128x128	+6	+0	0	-5	$2.5 \cdot 10^{-7}$	0.004	0.196
128x128x128	+6	+0	0	-6	$2.5 \cdot 10^{-7}$	0.024	0.356
128x128x128	+6	+0	0	-7	$2.5 \cdot 10^{-7}$	0.068	0.492
128x128x128	+6	+0	0	-8	$2.5 \cdot 10^{-7}$	0.108	0.656
128x128x128	+6	+0	0	-9	$2.5 \cdot 10^{-7}$	0.144	0.684
128x128x128	+6	+0	0	-10	$2.5 \cdot 10^{-7}$	0.144	0.788
128x128x128	+6	+0	0	-11	$2.5 \cdot 10^{-7}$	0.272	0.864
128x128x128	+6	+0	0	-12	$2.5 \cdot 10^{-7}$	0.324	0.912
128x128x128	+6	+0	0	-13	$2.5 \cdot 10^{-7}$	0.332	0.916
128x128x128	+6	+0	0	-14	$2.5 \cdot 10^{-7}$	0.472	0.940
128x128x128	+6	+0	0	-15	$2.5 \cdot 10^{-7}$	0.536	0.948
128x128x128	+8	+0	0	-5	$2.5 \cdot 10^{-7}$	0.000	0.024
128x128x128	+8	+0	0	-6	$2.5 \cdot 10^{-7}$	0.004	0.060
128x128x128	+8	+0	0	-7	$2.5 \cdot 10^{-7}$	0.028	0.056
128x128x128	+8	+0	0	-8	$2.5 \cdot 10^{-7}$	0.028	0.132
128x128x128	+8	+0	0	-9	$2.5 \cdot 10^{-7}$	0.028	0.200
128x128x128	+8	+0	0	-10	$2.5 \cdot 10^{-7}$	0.076	0.296
128x128x128	+8	+0	0	-11	$2.5 \cdot 10^{-7}$	0.084	0.336
128x128x128	+8	+0	0	-12	$2.5 \cdot 10^{-7}$	0.088	0.468
128x128x128	+8	+0	0	-13	$2.5 \cdot 10^{-7}$	0.084	0.544
128x128x128	+8	+0	0	-14	$2.5 \cdot 10^{-7}$	0.120	0.664
128x128x128	+8	+0	0	-15	$2.5 \cdot 10^{-7}$	0.136	0.700
128x128x128	+10	+0	0	-5	$2.5 \cdot 10^{-7}$	0.004	0.008
128x128x128	+10	+0	0	-6	$2.5 \cdot 10^{-7}$	0.000	0.020
128x128x128	+10	+0	0	-7	$2.5 \cdot 10^{-7}$	0.008	0.020
128x128x128	+10	+0	0	-8	$2.5 \cdot 10^{-7}$	0.016	0.036
128x128x128	+10	+0	0	-9	$2.5 \cdot 10^{-7}$	0.024	0.064
128x128x128	+10	+0	0	-10	$2.5 \cdot 10^{-7}$	0.028	0.096
128x128x128	+10	+0	0	-11	$2.5 \cdot 10^{-7}$	0.040	0.108
128x128x128	+10	+0	0	-12	$2.5 \cdot 10^{-7}$	0.032	0.140
128x128x128	+10	+0	0	-13	$2.5 \cdot 10^{-7}$	0.020	0.164
128x128x128	+10	+0	0	-14	$2.5 \cdot 10^{-7}$	0.036	0.208
128x128x128	+10	+0	0	-15	$2.5 \cdot 10^{-7}$	0.076	0.260

For reinforcement values  $\geq 16$ , the highest run success rates in the following table occur when the punishment is -24, the most extreme punishment in the table. One might do better in these cases with yet higher punishments.

game	r 1	r 2	p 1	p 2	plays/run	rsr 25 iters	rsr 50 iters	rsr 80 iters	rsr 100 iters
256x256x256	+11	+0	0	-10	$5 \cdot 10^{-7}$	0.004	0.360	0.264	0.264
256x256x256	+11	+0	0	-12	$5 \cdot 10^{-7}$	0.024	0.352	0.116	0.116
256x256x256	+11	+0	0	-14	$5 \cdot 10^{-7}$	0.052	0.308	0.040	0.040
256x256x256	+11	+0	0	-16	$5 \cdot 10^{-7}$	0.060	0.396	0.000	0.000
256x256x256	+11	+0	0	-18	$5 \cdot 10^{-7}$	0.124	0.344	0.000	0.000
256x256x256	+11	+0	0	-20	$5 \cdot 10^{-7}$	0.228	0.316	0.000	0.000
256x256x256	+11	+0	0	-22	$5 \cdot 10^{-7}$	0.264	0.236	0.000	0.000

game	r 1	r 2	p 1	p 2	plays/run	rsr 25 iters	rsr 50 iters	rsr 80 iters	rsr 100 iters
256x256x256	+11	+0	0	-24	$5 \cdot 10^{-7}$	0.284	0.060	0.000	0.000
256x256x256	+12	+0	0	-10	$5 \cdot 10^{-7}$	0.004	0.192	0.356	0.356
256x256x256	+12	+0	0	-12	$5 \cdot 10^{-7}$	0.008	0.308	0.236	0.236
256x256x256	+12	+0	0	-14	$5 \cdot 10^{-7}$	0.036	0.296	0.212	0.212
256x256x256	+12	+0	0	-16	$5 \cdot 10^{-7}$	0.040	0.384	0.204	0.204
256x256x256	+12	+0	0	-18	$5 \cdot 10^{-7}$	0.064	0.452	0.012	0.012
256x256x256	+12	+0	0	-20	$5 \cdot 10^{-7}$	0.084	0.396	0.000	0.000
256x256x256	+12	+0	0	-22	$5 \cdot 10^{-7}$	0.132	0.384	0.000	0.000
256x256x256	+12	+0	0	-24	$5 \cdot 10^{-7}$	0.252	0.348	0.000	0.000
256x256x256	+13	+0	0	-10	$5 \cdot 10^{-7}$	0.000	0.092	0.356	0.356
256x256x256	+13	+0	0	-12	$5 \cdot 10^{-7}$	0.008	0.148	0.404	0.404
256x256x256	+13	+0	0	-14	$5 \cdot 10^{-7}$	0.004	0.316	0.372	0.372
256x256x256	+13	+0	0	-16	$5 \cdot 10^{-7}$	0.016	0.340	0.300	0.300
256x256x256	+13	+0	0	-18	$5 \cdot 10^{-7}$	0.044	0.412	0.304	0.304
256x256x256	+13	+0	0	-20	$5 \cdot 10^{-7}$	0.076	0.412	0.236	0.236
256x256x256	+13	+0	0	-22	$5 \cdot 10^{-7}$	0.112	0.452	0.024	0.024
256x256x256	+13	+0	0	-24	$5 \cdot 10^{-7}$	0.112	0.464	0.000	0.000
256x256x256	+14	+0	0	-10	$5 \cdot 10^{-7}$	0.000	0.036	0.344	0.344
256x256x256	+14	+0	0	-12	$5 \cdot 10^{-7}$	0.004	0.084	0.400	0.400
256x256x256	+14	+0	0	-14	$5 \cdot 10^{-7}$	0.004	0.188	0.344	0.344
256x256x256	+14	+0	0	-16	$5 \cdot 10^{-7}$	0.008	0.252	0.360	0.360
256x256x256	+14	+0	0	-18	$5 \cdot 10^{-7}$	0.000	0.348	0.384	0.384
256x256x256	+14	+0	0	-20	$5 \cdot 10^{-7}$	0.016	0.372	0.332	0.332
256x256x256	+14	+0	0	-22	$5 \cdot 10^{-7}$	0.076	0.452	0.336	0.336
256x256x256	+14	+0	0	-24	$5 \cdot 10^{-7}$	0.044	0.484	0.292	0.292
256x256x256	+15	+0	0	-10	$5 \cdot 10^{-7}$	0.000	0.020	0.248	0.248
256x256x256	+15	+0	0	-12	$5 \cdot 10^{-7}$	0.000	0.048	0.316	0.316
256x256x256	+15	+0	0	-14	$5 \cdot 10^{-7}$	0.000	0.108	0.412	0.412
256x256x256	+15	+0	0	-16	$5 \cdot 10^{-7}$	0.004	0.156	0.480	0.480
256x256x256	+15	+0	0	-18	$5 \cdot 10^{-7}$	0.020	0.268	0.428	0.428
256x256x256	+15	+0	0	-20	$5 \cdot 10^{-7}$	0.032	0.336	0.420	0.420
256x256x256	+15	+0	0	-22	$5 \cdot 10^{-7}$	0.024	0.424	0.428	0.428
256x256x256	+15	+0	0	-24	$5 \cdot 10^{-7}$	0.048	0.452	0.364	0.364
256x256x256	+16	+0	0	-10	$5 \cdot 10^{-7}$	0.004	0.008	0.112	0.112
256x256x256	+16	+0	0	-12	$5 \cdot 10^{-7}$	0.000	0.036	0.264	0.264
256x256x256	+16	+0	0	-14	$5 \cdot 10^{-7}$	0.012	0.032	0.388	0.388
256x256x256	+16	+0	0	-16	$5 \cdot 10^{-7}$	0.004	0.112	0.416	0.416
256x256x256	+16	+0	0	-18	$5 \cdot 10^{-7}$	0.004	0.152	0.448	0.448
256x256x256	+16	+0	0	-20	$5 \cdot 10^{-7}$	0.020	0.256	0.432	0.432
256x256x256	+16	+0	0	-22	$5 \cdot 10^{-7}$	0.016	0.320	0.444	0.444
256x256x256	+16	+0	0	-24	$5 \cdot 10^{-7}$	0.036	0.352	0.532	0.532
256x256x256	+17	+0	0	-10	$5 \cdot 10^{-7}$	0.000	0.012	0.076	0.076
256x256x256	+17	+0	0	-12	$5 \cdot 10^{-7}$	0.000	0.008	0.192	0.192
256x256x256	+17	+0	0	-14	$5 \cdot 10^{-7}$	0.000	0.036	0.324	0.324
256x256x256	+17	+0	0	-16	$5 \cdot 10^{-7}$	0.008	0.040	0.320	0.320
256x256x256	+17	+0	0	-18	$5 \cdot 10^{-7}$	0.008	0.140	0.456	0.456
256x256x256	+17	+0	0	-20	$5 \cdot 10^{-7}$	0.008	0.176	0.492	0.492
256x256x256	+17	+0	0	-22	$5 \cdot 10^{-7}$	0.024	0.224	0.472	0.472
256x256x256	+17	+0	0	-24	$5 \cdot 10^{-7}$	0.036	0.320	0.508	0.508



game	r 1	r 2	p 1	p 2	plays/run	rsr 25 iters	rsr 50 iters	rsr 80 iters	rsr 100 iters
256x256x256	+18	+0	0	-10	$5 \cdot 10^7$	0.000	0.004	0.088	0.088
256x256x256	+18	+0	0	-12	$5 \cdot 10^7$	0.000	0.024	0.132	0.132
256x256x256	+18	+0	0	-14	$5 \cdot 10^7$	0.004	0.028	0.232	0.232
256x256x256	+18	+0	0	-16	$5 \cdot 10^7$	0.000	0.028	0.356	0.356
256x256x256	+18	+0	0	-18	$5 \cdot 10^7$	0.004	0.100	0.408	0.408
256x256x256	+18	+0	0	-20	$5 \cdot 10^7$	0.008	0.128	0.468	0.468
256x256x256	+18	+0	0	-22	$5 \cdot 10^7$	0.004	0.160	0.476	0.476
256x256x256	+18	+0	0	-24	$5 \cdot 10^7$	0.016	0.220	0.484	0.484
256x256x256	+19	+0	0	-10	$5 \cdot 10^7$	0.000	0.008	0.044	0.044
256x256x256	+19	+0	0	-12	$5 \cdot 10^7$	0.000	0.004	0.100	0.100
256x256x256	+19	+0	0	-14	$5 \cdot 10^7$	0.000	0.016	0.152	0.152
256x256x256	+19	+0	0	-16	$5 \cdot 10^7$	0.008	0.020	0.248	0.248
256x256x256	+19	+0	0	-18	$5 \cdot 10^7$	0.008	0.044	0.352	0.352
256x256x256	+19	+0	0	-20	$5 \cdot 10^7$	0.004	0.060	0.408	0.408
256x256x256	+19	+0	0	-22	$5 \cdot 10^7$	0.016	0.108	0.436	0.436
256x256x256	+19	+0	0	-24	$5 \cdot 10^7$	0.020	0.132	0.456	0.456
256x256x256	+20	+0	0	-10	$5 \cdot 10^7$	0.000	0.000	0.024	0.024
256x256x256	+20	+0	0	-12	$5 \cdot 10^7$	0.000	0.012	0.076	0.076
256x256x256	+20	+0	0	-14	$5 \cdot 10^7$	0.004	0.004	0.104	0.104
256x256x256	+20	+0	0	-16	$5 \cdot 10^7$	0.000	0.028	0.188	0.188
256x256x256	+20	+0	0	-18	$5 \cdot 10^7$	0.000	0.064	0.308	0.308
256x256x256	+20	+0	0	-20	$5 \cdot 10^7$	0.008	0.052	0.416	0.416
256x256x256	+20	+0	0	-22	$5 \cdot 10^7$	0.020	0.088	0.328	0.328
256x256x256	+20	+0	0	-24	$5 \cdot 10^7$	0.016	0.148	0.480	0.480

The following table is for a hybrid dynamics where  $r1 = r2$ , but  $p1 \neq p2$ .

game	r 1	r 2	p 1	p 2	plays/run	rsr 16 iters
32x32x32	+1	+1	0	0	$10^7$	0.025
32x32x32	+1	+1	0	-1	$10^7$	0.181
32x32x32	+1	+1	0	-2	$10^7$	0.288
32x32x32	+1	+1	0	-3	$10^7$	0.443
32x32x32	+1	+1	0	-4	$10^7$	0.533
32x32x32	+1	+1	0	-5	$10^7$	0.616
32x32x32	+1	+1	-1	0	$10^7$	0.183
32x32x32	+1	+1	-1	-1	$10^7$	0.000
32x32x32	+1	+1	-1	-2	$10^7$	0.000
32x32x32	+1	+1	-1	-3	$10^7$	0.000
32x32x32	+1	+1	-1	-4	$10^7$	0.000
32x32x32	+1	+1	-1	-5	$10^7$	0.000
32x32x32	+1	+1	-2	0	$10^7$	0.286
32x32x32	+1	+1	-2	-1	$10^7$	0.000
32x32x32	+1	+1	-2	-2	$10^7$	0.000
32x32x32	+1	+1	-2	-3	$10^7$	0.000
32x32x32	+1	+1	-2	-4	$10^7$	0.000
32x32x32	+1	+1	-2	-5	$10^7$	0.000
32x32x32	+1	+1	-3	0	$10^7$	0.438
32x32x32	+1	+1	-3	-1	$10^7$	0.000

game	r 1	r 2	p 1	p 2	plays/run	rsr 16 iters
32x32x32	+1	+1	-3	-2	10 <sup>7</sup>	0.000
32x32x32	+1	+1	-3	-3	10 <sup>7</sup>	0.000
32x32x32	+1	+1	-3	-4	10 <sup>7</sup>	0.000
32x32x32	+1	+1	-3	-5	10 <sup>7</sup>	0.000
32x32x32	+2	+2	0	0	10 <sup>7</sup>	0.033
32x32x32	+2	+2	0	-1	10 <sup>7</sup>	0.109
32x32x32	+2	+2	0	-2	10 <sup>7</sup>	0.169
32x32x32	+2	+2	0	-3	10 <sup>7</sup>	0.204
32x32x32	+2	+2	0	-4	10 <sup>7</sup>	0.280
32x32x32	+2	+2	0	-5	10 <sup>7</sup>	0.322
32x32x32	+2	+2	-1	0	10 <sup>7</sup>	0.097
32x32x32	+2	+2	-1	-1	10 <sup>7</sup>	0.009
32x32x32	+2	+2	-1	-2	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-1	-3	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-1	-4	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-1	-5	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-2	0	10 <sup>7</sup>	0.173
32x32x32	+2	+2	-2	-1	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-2	-2	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-2	-3	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-2	-4	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-2	-5	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-3	0	10 <sup>7</sup>	0.211
32x32x32	+2	+2	-3	-1	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-3	-2	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-3	-3	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-3	-4	10 <sup>7</sup>	0.000
32x32x32	+2	+2	-3	-5	10 <sup>7</sup>	0.000
32x32x32	+3	+3	0	0	10 <sup>7</sup>	0.035
32x32x32	+3	+3	0	-1	10 <sup>7</sup>	0.071
32x32x32	+3	+3	0	-2	10 <sup>7</sup>	0.114
32x32x32	+3	+3	0	-3	10 <sup>7</sup>	0.181
32x32x32	+3	+3	0	-4	10 <sup>7</sup>	0.195
32x32x32	+3	+3	0	-5	10 <sup>7</sup>	0.231
32x32x32	+3	+3	-1	0	10 <sup>7</sup>	0.913
32x32x32	+3	+3	-1	-1	10 <sup>7</sup>	0.922
32x32x32	+3	+3	-1	-2	10 <sup>7</sup>	0.909
32x32x32	+3	+3	-1	-3	10 <sup>7</sup>	0.926
32x32x32	+3	+3	-1	-4	10 <sup>7</sup>	0.902
32x32x32	+3	+3	-1	-5	10 <sup>7</sup>	0.922
32x32x32	+3	+3	-2	0	10 <sup>7</sup>	0.137
32x32x32	+3	+3	-2	-1	10 <sup>7</sup>	0.901
32x32x32	+3	+3	-2	-2	10 <sup>7</sup>	0.221

game	r 1	r 2	p 1	p 2	plays/run	rsr 25 iters
64x64x64	+1	+1	0	0	2x10 <sup>7</sup>	0.000
64x64x64	+1	+1	0	-1	2x10 <sup>7</sup>	0.045

game	r 1	r 2	p 1	p 2	plays/run	rsr 25 iters
64x64x64	+1	+1	0	-2	2x10 <sup>7</sup>	0.330
64x64x64	+1	+1	0	-3	2x10 <sup>7</sup>	0.565
64x64x64	+1	+1	0	-4	2x10 <sup>7</sup>	0.660
64x64x64	+1	+1	0	-5	2x10 <sup>7</sup>	0.640
64x64x64	+1	+1	-1	0	2x10 <sup>7</sup>	0.050
64x64x64	+1	+1	-1	-1	2x10 <sup>7</sup>	0.000
64x64x64	+1	+1	-1	-2	2x10 <sup>7</sup>	0.000
64x64x64	+1	+1	-1	-3	2x10 <sup>7</sup>	0.000
64x64x64	+1	+1	-1	-4	2x10 <sup>7</sup>	0.000
64x64x64	+1	+1	-1	-5	2x10 <sup>7</sup>	0.000
64x64x64	+2	+2	0	0	2x10 <sup>7</sup>	0.000
64x64x64	+2	+2	0	-1	2x10 <sup>7</sup>	0.000
64x64x64	+2	+2	0	-2	2x10 <sup>7</sup>	0.010
64x64x64	+2	+2	0	-3	2x10 <sup>7</sup>	0.035
64x64x64	+2	+2	0	-4	2x10 <sup>7</sup>	0.040
64x64x64	+2	+2	0	-5	2x10 <sup>7</sup>	0.070
64x64x64	+2	+2	-1	0	2x10 <sup>7</sup>	0.000
64x64x64	+2	+2	-1	-1	2x10 <sup>7</sup>	0.000
64x64x64	+2	+2	-1	-2	2x10 <sup>7</sup>	0.000
64x64x64	+2	+2	-1	-3	2x10 <sup>7</sup>	0.000
64x64x64	+2	+2	-1	-4	2x10 <sup>7</sup>	0.000
64x64x64	+2	+2	-1	-5	2x10 <sup>7</sup>	0.000

## REFERENCES

- J. M. Alexander, B. Skyrms, S. L. Zabell (2012). Inventing new signals. *Dynamic Games and Applications*, 2, 129–145.
- Argiento, Raffaele, Robin Pemantle, Brian Skyrms and Stas Volkov (2009). Learning to signal: analysis of a micro-level reinforcement model. *Stochastic Processes and their Applications*, 119(2), 373–390
- Barrett, Jeffrey A. (2007a). Dynamic partitioning and the conventionality of kinds. *Philosophy of Science*, 74, 527–546.
- Barrett, Jeffrey A. (2007b). The evolution of coding in signaling games. *Theory and Decision*, 67(2), 223–237.
- Barrett, Jeffrey A. (2006). Numerical simulations of the Lewis signaling game: learning strategies, pooling equilibria, and the evolution of grammar. *UC Irvine Institute for Mathematical Behavioral Sciences Technical Report*.  
<https://www.imbs.uci.edu/research/technical.php>.
- Barrett, Jeffrey A and Kevin Zollman (2009). The role of forgetting in the evolution and learning of language. *Journal of Experimental and Theoretical Artificial Intelligence*, 21(4), 293–309.
- Barrett, J. A., C. T. Cochran, S. Huttegger, N. Fujiwara (2017) Hybrid learning in signaling games. *Journal of Experimental & Theoretical Artificial Intelligence* 29(5), 1–9.
- Cochran, C. T. and J. A. Barrett (2022) The Efficacy of Human Learning in Lewis-Skyrms Signaling Games. forthcoming in *British Journal for the Philosophy of Science*.
- Cochran, C. T. and J. A. Barrett (2021) How Signaling Conventions are Established. *Synthese*, 199, 4367–4391.
- Erev, I. and A. E. Roth (1998) Predicting how people play games: reinforcement learning in experimental games with unique, mixed strategy equilibria. *American Economic Review*, 88, 848–81.
- Herrnstein, R. J. (1970) On the law of effect. *Journal of the Experimental Analysis of Behavior*, 13, 243–266.
- Hofbauer, J. and S. Huttegger (2008) Feasibility of communication in binary signaling games. *Journal of Theoretical Biology*, 254(4), 843–849.
- Lewis, David (1969). *Convention*. Cambridge, MA: Harvard University Press.
- Roth, A. E. and I. Erev (1995) Learning in extensive form games: experimental data and simple dynamical models in the intermediate term. *Games and Economic Behavior*, 8, 164–212.
- Skyrms, Brian (2010). *Signals: Evolution, Learning, & Information*, New York: Oxford University Press.
- Skyrms, Brian (2006). Signals. *Philosophy of Science*, 75(5), 489–500.