

# The Effects of Trust on Convention Emergence

Nir Oren<sup>1</sup>, Nathan Griffiths<sup>2</sup>, and Michael Luck<sup>3</sup>

<sup>1</sup> Department of Computing Science, University of Aberdeen  
Aberdeen, AB24 3UE, Scotland  
n.oren@abdn.ac.uk

<sup>2</sup> Department of Computer Science, University of Warwick  
Coventry, CV4 7AL, UK  
nathan@dcs.warwick.ac.uk

<sup>3</sup> Department of Informatics, King's College London  
London, WC2R 2LS, UK  
michael.luck@kcl.ac.uk

**Abstract.** In this paper, we consider the effects of trust and conventions on the evolution of a multi-agent system. To do so, we provide a simulation in which agents are able to use a trust mechanism to select interaction partners, and a learning mechanism that leads to convention emergence. We examine the impacts on the system of: different network topologies; the presence or absence of malicious agents; and the presence or absence of the trust and learning system. Our results indicate that while trust has a slight positive impact on the rate of convergence emergence, its main benefit arises when malicious agents are present in the system.

## 1 Introduction

Trust, norms and conventions all aim to increase the likelihood of globally desirable system behaviour. The former operates by reducing the probability that a poorly behaving agent will continue to be selected for interaction by others, while both norms and conventions aim to specify correct patterns of behaviour, with norms imposing sanctions on agents who deviate from this behaviour.

Within human societies, behaviour is regulated using a mix of these mechanisms. However since any mechanism typically has some cost associated with it (e.g. requiring monitoring and sanctioning in the case of norms, or potentially disadvantaging new agents in the case of trust), the question arises of why several distinct behaviour regulation mechanisms have emerged, and whether one is sufficient to maintain desirable system outcomes. This is important, because if we are able to achieve desired behaviour through only a subset of these mechanisms, then in principle we may be able to do so more efficiently.

The impact of trust on a system in which convention emergence occurs under different system properties is therefore crucial to understand clearly. That is the focus of this paper. More specifically, we aim to establish the impacts of different network topologies on agent outcomes in such systems, as well the effects of different types of malicious agents and of *churn* — the replacement of existing agents by new agents. We investigate the following core hypotheses.

1. The topology of the network can affect the usefulness of a trust mechanism.
2. Trust can cause islands to form, inhibiting the emergence of conventions.
3. Trust is critical for correct system behaviour when malicious agents exist.

We seek to validate these hypotheses through simulation in which a population of agents interact with each other playing a coordination game. A trust mechanism enables agents to select partners to interact with, while a learning mechanism allows agents to learn a strategy, simulating the emergence of conventions<sup>1</sup>. We evaluate the effects of the presence or absence of each of these mechanisms on different configurations of the system, considering different interaction network topologies, malicious agents, and churn.

Our results validate our hypotheses, indicating that the main impact of trust arises when malicious agents are present, but that its use can also hamper global convergence emergence. In the following section we introduce the notions of trust and conventions, and identify the key related work. Section 3 describes the cooperative game in which agents participate. Section 4.1 introduces our simulation system in more detail, and our experimental results are described in Section 4.2. Finally, Section 5 concludes the paper.

## 2 Trust and Conventions

In open dynamic decentralised systems, global control is not always practical, and agents need to decide for themselves who to interact with, and what behaviour to exhibit. In multi-agent systems (MAS), the problem of *who* to interact with has often been addressed using trust and reputation to build an estimation of how others are likely to act, and so enable trustworthy partners to be selected. In many domains, answering the question of *what* action to perform requires consideration of the actions that other individuals are likely to take. For example, consider deciding whether to drive on the left or right side of the road on arrival in an unfamiliar location. There is no intrinsic preference between the options, but there is clearly a very strong incentive to drive on the same side as others to prevent collisions. This problem of selecting between actions, when their effectiveness is strongly influenced by the choices of others, is addressed by the notion of conventions, where a convention can be viewed as a social rule or standard of behaviour agreed on, or adopted by, a set of individuals [7, 22].

There is a significant body of work on trust and reputation in multi-agent systems, ranging from lightweight image scoring approaches to rich models using detailed histories of previous transactions. Conventions have also been explored in relation to multi-agent systems, typically considering how they can emerge as a result of individuals learning which actions afford them the highest utility. However, to our knowledge, there has been little investigation of the relationship between trust and conventions. In the remainder of this section we briefly introduce the key literature on trust and conventions, as related to this paper.

---

<sup>1</sup> In this paper we focus on conventions as they are a broader notion than norms, which typically require sanctions or rewards.

## 2.1 Trust in Multi-Agent Systems

Trust and reputation are effective mechanisms for supporting cooperative and coordinated behaviour by enabling agents to choose who to interact with [6, 12, 14, 15]. We can view trust as a subjectively held belief about the likelihood with which another agent will act as expected, typically meaning that the other individual will be cooperative if given an opportunity to defect and receive higher payoff. Reputation can be viewed as a socially accepted trust assessment for a given individual [6]. Trust often requires significant historical interaction data for accurate assessments, and instantiations often make use of multiple dimensions of information [4, 16]. Reputation is a key mechanism when individuals have insufficient direct interaction history with which to assess a partner's trust value, but it still relies on being able to draw on the indirect experiences of others.

The trust and reputation mechanisms that have exhibited the most promising results tend to be the most complex, requiring significant information on agents' previous interactions [4, 16]. In decentralised open multi-agent systems, especially where individuals can leave and join the system, such mechanisms may be less suitable than alternatives with less complex requirements. Indeed, it has been shown that very simple trust mechanisms with low overheads can still be effective [12, 14].

Nowak and Sigmund introduced image scoring as a simple instantiation of trust and reputation with low overheads, based on the notion of indirect reciprocity, in which cooperation emerges without requiring subsequent interactions between the same individuals [12, 13]. This property is key to its suitability in open decentralised systems. Each agent maintains an image score for each individual it interacts with or observes interacting. Cooperative actions increase the image score by one, and selfish actions decrease it by one. When deciding whether to cooperate or not, an agent compares its strategy, an integer, with the perceived image score of the potential partner (if no data is available, it is assumed to be zero). If the image score is greater than or equal to the strategy then the agent cooperates. We use this image scoring view of trust as the basis for the trust model used in this paper.

## 2.2 Conventions in Multi-Agent Systems

A convention is a social rule or standard of behaviour agreed upon by a set of individuals [7, 22], and can be considered established once a high proportion of a population adheres to it for a significant amount of the time [7]. Conventions can increase levels of coordination in multi-agent systems [2, 5, 21], and they are a powerful abstraction tool for modelling the aggregate interactions of agents.

Conventions can be generated offline by system designers or dynamically emerge through interactions. However, offline generation is often impractical due to limited knowledge of society characteristics, time variance, and computational expense. Moreover, such conventions also lack robustness and, as a result, much research has concentrated on generating conventions online [11, 17].

A common theme in the definition of conventions revolves around the regularity in the behaviour of a population in repeated iterations in the same situation. From a multi-agent systems perspective, conventions have been viewed in game-theoretic terms, as a

		Agent 1	
	action	A	B
Agent 2	A	1	-1
	B	-1	1

(a)

		Agent 1	
	action	A	B
Agent 2	A	-1	1
	B	1	-1

(b)

**Table 1.** Payoff matrices for (a) normal agents and (b) imperfect malicious agents.

restriction of agents’ decisions to a single choice in a coordination game [19], or more generally being defined by a high proportion of agents adopting a particular strategy [7].

Convention emergence is often illustrated by considering coordination games, such as which side of the road to drive on [11, 18]. The ideal is for every agent to adhere to the same convention, but this may be unrealistic. As the multitude of global conventions in real-world traffic systems show, it is not necessary for a single global convention to pervade for high levels of local coordination to emerge. However, the cost of inappropriate or inefficient conventions is very high.

It is well known that constraints on interactions between agents, originating from different network topologies, impact on how conventions emerge. For example, Villatoro *et al.* have shown that fully connected networks are the quickest to converge, with scale-free being significantly slower [20]. Other researchers have also shown that topology is an important factor in convention emergence, but have often assumed full observability of others [2, 7].

### 3 The Cooperation Game

In order to be able to investigate the impact of trust in the context of convention emergence, we instantiate a simple MAS in which agents interact with each other by playing a cooperative two party coordination game<sup>2</sup>; this is equivalent to the setting considered in [18]. Agents within this system are described through a single parameter — the probability of playing an action within an interaction (i.e. the likelihood of selecting a specific action within the game). Each agent is able to interact with a subset of the other agents in the system, which we refer to as an agent’s *neighbourhood*. We encode this structure through an undirected *interaction network* or graph, in which agents are represented as nodes, and the ability to interact is represented by an edge.

The game progresses in discrete rounds. During each round, every agent selects others in its neighbourhood to interact with. Each interaction consists of both agents selecting an action. These actions are evaluated through a payoff matrix (shown in Table 1(a)), and the agents are then informed of their reward.

Clearly, making the right choice of which neighbours to interact with is the crucial aspect of this game, where some may choose compatible actions, and others incompatible actions. In this game, different methods exist for neighbour selection (e.g. weighting

<sup>2</sup> Note that there is no intrinsic preference between actions, differentiating this work from the evolution of trust for prisoner’s dilemma games [9].

---

**Algorithm 1** The PHC Q-learning algorithm.

---

```
1: function Q-LEARNING( $\epsilon, \alpha, \delta$ )
2:    $Q[A], Q[B] \leftarrow 0$ 
3:    $\pi \leftarrow$  uniform random number in  $[0,1]$ 
4:   while true do
5:     if a uniform random number selected from  $[0,1] > \epsilon$  then
6:       select action  $A$  with probability  $\pi$ , otherwise select  $B$ 
7:     else
8:       select action from  $A, B$  with uniform likelihood
9:     end if
10:     $r \leftarrow$  reward from the interaction
11:     $Q[\text{selected action}] = (1 - \alpha)Q[\text{selected action}] + \alpha r$ 
12:    if  $Q[A] > Q[B]$  then  $\pi \leftarrow \pi + \delta$ 
13:    if  $Q[B] > Q[A]$  then  $\pi \leftarrow \pi - \delta$ 
14:    ensure that  $\pi \in [0, 1]$ 
15:  end while
16: end function
```

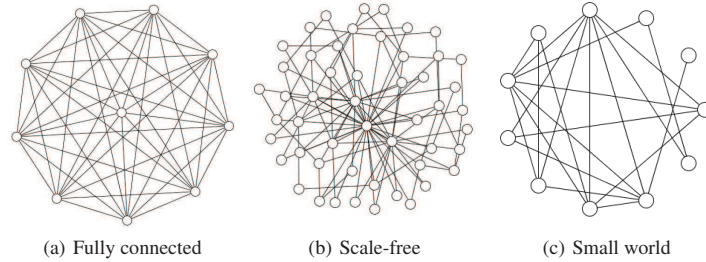
---

this selection based on a trust mechanism, or randomly choosing a neighbour for interaction), and for updating the probability of playing a specific action (e.g. keeping this probability fixed, or performing learning to modify it). Moreover, the game may be played with different interaction topologies, different numbers of malicious agents, and agents leaving and entering the system, which we refer to as *churn*. Next, we describe each of these aspects.

**Trust Mechanism** Our trust mechanism stores the outcome of the last  $n$  interactions. Given a total of  $s$  successful interactions with all neighbours, the likelihood of interacting with a specific neighbour is  $x/s$  where  $x$  is the number of successful interactions with that neighbour, or  $x = 1$ , whichever is greater. This is a simple trust mechanism, and although more sophisticated mechanisms are possible, they typically have comparatively large overheads in terms of information and computation requirements.

**Convention Learning** After every interaction, an agent can update its probability of playing an action. To do so, the agent utilises a slightly modified version of 2 action PHC Q-learning [1] with no lookahead, as described in Algorithm 1, and as has previously been used when investigating conventions [18]. Unlike standard PHC Q-learning, we randomise the initial policy.

**Malicious Agents and Churn** Malicious agents seek to undermine the system, either due to holding a different utility function to the other agents, or out of malice. In this paper, we consider two types of malicious agents, referred to as *imperfect* or *omniscient* malicious agents respectively. The former type of agent utilises the utility function shown in Table 1(b), and otherwise acts as an agent able to utilise both trust and convention learning. Omniscient malicious agents always cause the agent interacting with them to obtain a utility of -1. While it is difficult to identify an obvious real



**Fig. 1.** Illustrative network topologies

world analogy to the latter class of agent, they represent a pathological worst case, and are thus useful in our evaluation.

Churn represents the departure of existing agents, and the introduction of new agents as the system runs. These new agents are unaware of any conventions that may have emerged, and also disrupt the trust mechanism. We simulate churn by randomly selecting a small proportion of the population and resetting their interaction histories and learned parameters.

**Network Topologies** While many different network structures are possible, we focus on three archetypes of interaction topologies, namely fully connected networks, scale-free networks, and small world networks. Fully connected networks provide agents with the greatest latitude in choosing interaction partners, allowing them to easily ignore malicious agents through the use of a trust mechanism. However, such networks are, in most situations unrealistic, and we therefore consider additional topologies, which model the properties of a large family of real world systems such as social networks, citation networks, and road networks. Small world networks are constructed from a ring topology. A parameter  $k$  adds extra edges to ensure that all agents within  $k/2$  hops from each other are connected. Another parameter, the rewiring probability  $\beta$ , is then used to replace some of these connections with others to randomly selected nodes in the network [8]. The connections between nodes in scale-free networks follow the power law distribution so that some nodes have very many connections, but the majority have very few [3]. These networks are illustrated in Figure 1.

## 4 Evaluation

We begin this section by describing our simulation environment, following which we detail the experiments we ran, and discuss our results.

#### 4.1 The Simulation Environment

To investigate the impact of these mechanisms, we developed a simulation environment for the coordination game, and undertook experiments consisting of running repeated iterations within this environment. During each iteration, every agent in the system selects another agent for a single interaction.

The trust mechanism was instantiated with a history of 20 interactions ( $n = 20$ ), so that agents had this number of interactions with others available for use when determining trust, and our convention learning approach specified in Algorithm 1 was instantiated with the exploration parameter  $\epsilon$  being set to 0.02, while the  $\alpha$  and  $\delta$  learning parameters were set to 0.1 and 0.01 respectively. The full set of parameters can be found in Table 2. The simulation environment, and thus the experiments, consisted of 50 agents, averaged over 10 runs, with 10000 iterations per run.<sup>3</sup> When examining the effects of malicious agents and churn, we assigned each individual agent a probability (0.1) of becoming a malicious agent or being reset (to simulate churn), every 1000 rounds. Once an agent becomes malicious, it remains so until the simulation terminates.

Number of Agents	50	Number of iterations	10000
Number of runs	10	Exploration probability $\epsilon$	0.02
Likelihood of agent replacement	0.1	Frequency of agent replacement	1000
Trust system memory $n$	20	Q learning $\delta, \alpha$	0.01, 0.1
Scale-free graph $m$	2	Small world graph $\beta, k$	0.2, 6

Table 2. Simulation parameters

#### 4.2 Experiments

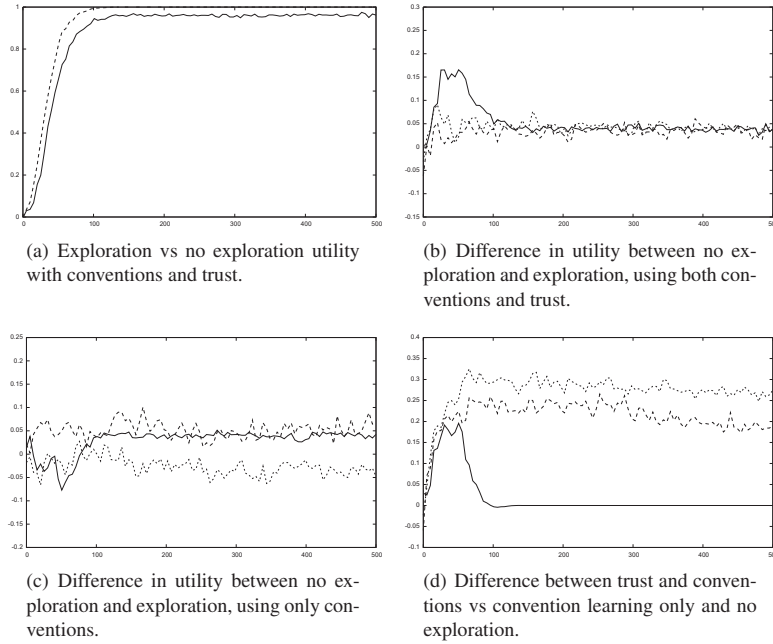
In this section, we provide further details of our experiments and results. Unless stated, results were obtained by averaging together 10 runs of the system.<sup>4</sup>

**Experiment 1: Topologies** We began by undertaking a high level comparison of the three topologies, comparing average utility where agents utilise both convention learning and trust, and where only the former is utilised, to evaluate the effect of trust. Figure 2(a) illustrates the utility obtained when exploration is present as part of the convention learning (solid), and when no exploration learning occurs, i.e.  $\epsilon = 0$  (the dashed line). Note that while utility converges to 1 for fully connected networks in the no-exploration case, this does not occur for any other network topologies. We believe that the core reason for this is the so-called *island effect*, discussed later.

More generally, Figure 2(b) illustrates the *difference* between the no-exploration and exploration utility curves for all network topologies. Given these curves, it appears

<sup>3</sup> Every agent will thus have at least 200 interactions in the course of a single run.

<sup>4</sup> For clarity, all graphs were smoothed using GNUPlot’s acspline option (weight=0.1).

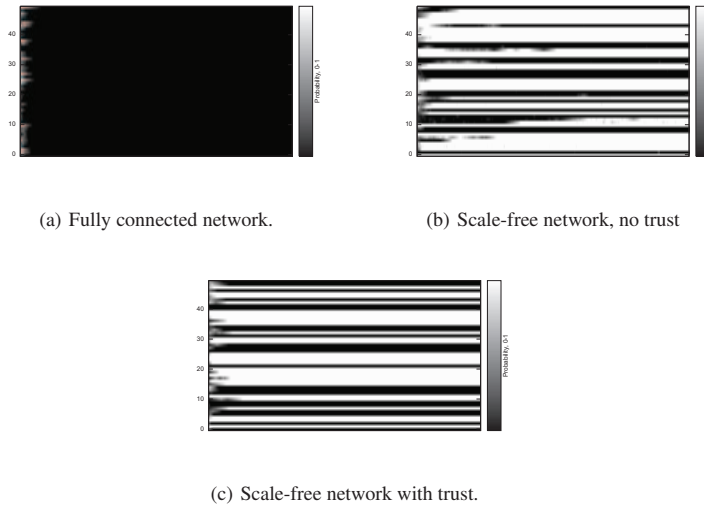


**Fig. 2.** Comparing average utilities with no malicious agents. In (b)–(d), the solid line indicates results for the fully connected graph, the dashed line indicates the small world network, and the dotted line the scale-free network.

that when no malicious agents are present, exploration is unnecessary; the average utility obtained with no exploration appears to exceed (or in the case of scale-free graphs, approximately equal) the utility obtained when exploration takes place. Furthermore, as demonstrated for example in Figure 2(a), the rate of convergence, particularly in the fully connected network, is faster when no exploration takes place. Figure 2(c) illustrates the difference in utilities when exploration does and does not occur, when only conventions (rather than both conventions and trust) are present in the system. Again, the exploration seems to detract from the rate at which conventions emerge. Given this, in the remainder of our experiments, we considered only the no-exploration form of Q-learning. It should be noted that only 500 iterations are displayed in most of our figures since that is sufficient to identify dominant system trends, though we ran up to 10000 iterations in our experiments.

Figure 2(d) illustrates the difference between utility curves when only conventions are used, and when both conventions and trust are used in the no-exploration scenario. While the difference for fully connected graphs quickly reaches zero, indicating fast



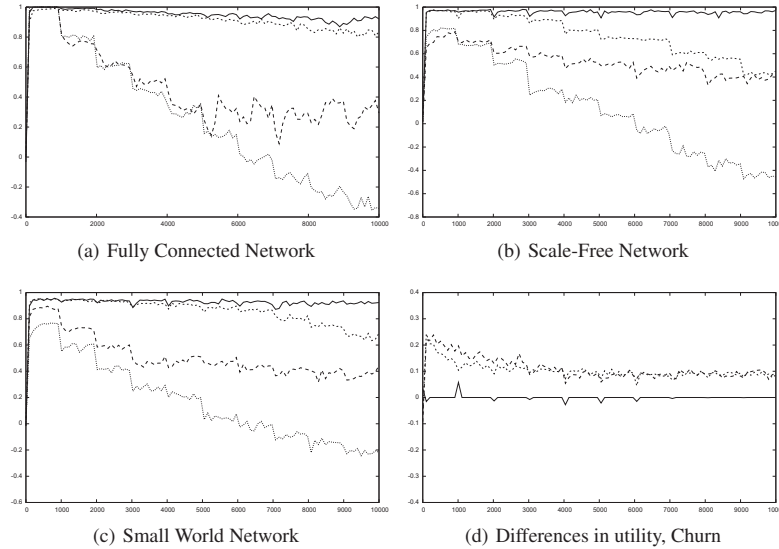


**Fig. 3.** Heat maps. The vertical axis identifies different agents in the system, while the horizontal axis is the iteration number. Plots are shown for the first 1000 iterations.

convergence in both the trust and no trust cases, this does not occur for other topologies. For example, scale-free networks neared convergence after approximately 10000 iterations.

We hypothesised that the presence of trust will lead to the formation of islands, where highly trusted cliques of agents form, causing all these agents to play a specific action, while other cliques play some other action. Figure 3(a) shows a heat map of the probabilities that agents play some move for the fully connected graph. It is clear that no islands form. This can be contrasted with Figure 3(b) which shows a heat map for the scale-free network case *where no trust mechanism exists*. Here, it is clear that distinct islands form. This suggests that islands form due to physical neighbourhoods rather than due to the trust mechanism. A closer examination of the trust system explains why.

Consider the fully connected case for a system of 50 agents. Initially, there is a 0.02 chance ( $1/50$ ) of an agent being selected for an interaction. Now consider the case where an agent was selected 20 times for interactions (the maximum possible due to our trust system's memory). In this case, there is approximately a 0.3 likelihood of this agent being selected for interaction, but there is a 0.7 likelihood that some other agent will be selected. As an agent's neighbourhood shrinks, this effect (i.e. that trust plays a minority role in interaction partner selection) becomes less pronounced, but the constraints due to lack of connectivity begin to dominate. Comparing Figure 3(b) and



**Fig. 4.** Average utility under difference scenarios. In (a)–(c) solid lines plot imperfect malicious agents using trust and conventions; long dashes plot imperfect malicious agents with no trust; short dashes plot omniscient malicious agents using trust and conventions; dots plot omniscient malicious agents with no trust. In (d), solid lines plot fully connected; long dashes plot small world; short dashes plot scale-free.

3(c), we see that the main effect of trust in the scale-free case arises from the increased rate of convergence, and that islands (represented as horizontal lines of one colour) appear in both cases. Additional work is needed to verify whether the island effect is greater in the presence of trust. We believe that this interplay between trust and limited interaction partners also explains the lack of convergence encountered in Figure 2(d).

**Experiment 2: Maliciousness and Churn** We have shown that in many situations, the use of trust in the context of convention learning either provides only limited help or perhaps even hinders the effectiveness of a normally functioning system. However, trust is aimed at situations where agents are not always benign. In this experiment, we examined the effectiveness of the use of trust in the presence of both omniscient and imperfect malicious agents.

Figures 4(a)–4(c) illustrate the effects of malicious agents on all network topologies in the presence and absence of trust. These can be contrasted with Figure 4(d), which was plotted for a system containing churn. While the different x-axis scale between this figure and 2(d) — required in the latter to identify relevant results — make it difficult to

see, the curves here are similar to those of the latter figure. Churn introduces discontinuities into the curve, which the trust and convention learning system is able to quickly overcome. Therefore while useful, trust has only a small effect in the presence of churn. However, its presence in situations where malicious agents exist is critical to the system. Note that due to our agent replacement strategy, approximately 65% of the agents were malicious after 9000 iterations. Also note that graph topologies had a relatively minor effect on the overall impact of malicious agents.

## 5 Discussion and Conclusions

With regard to our hypotheses, our experiments confirm that network topology can have an effect on the usefulness of a trust mechanism, and that trust is indeed critical in helping convergence emerge in the presence of malicious agents. Our experiments also indicate that, at least in the manner in which we implemented trust, the problem of islands of conventions is minor, with the dominant factor regarding the emergence of these islands being the topology of the network itself.

We have also shown that trust increases the initial rate of convergence of a convention. While Q-learning in our domain meant that conventions emerged relatively quickly, this property could be useful in some domains.

We intend to pursue several avenues of future work. First, we have only highlighted the most significant results of our work for very specific parameter settings. Exploring how other parameters affect the system could allow us to identify additional situations where trust is, or is not, necessary. In conjunction with this, we intend to extend our work to the normative domain. Following work such as [10], a norm requires an agent to adhere to some specific behaviour. However, unlike a convention, the violation of a norm allows for a sanction to be imposed on the violator. We intend to investigate why this sanctioning mechanism, whose effects are similar to that of a trust mechanism, emerged, and to investigate the interplay between these.

To our knowledge, our work is the first to seek to identify the relationship between conventions and trust, determining the situations in which one, or both, of these mechanisms is required to obtain desired system behaviour. We believe that several important avenues of future research remain, and that our work will aid in the design of efficient large scale multi-agent systems.

## References

1. M. Bowling and M. Veloso. Multiagent learning using a variable learning rate. *Artificial Intelligence*, 136(2):215–250, April 2002.
2. J. Delgado. Emergence of social conventions in complex networks. *Artificial Intelligence*, 141(1-2):171–185, October 2002.
3. D. Eppstein and J. Wang. A steady state model for graph power laws. *2nd Int. Workshop on Web Dynamics*, 2002.
4. D. Huynh, N. R. Jennings, and N. R. Shadbolt. An integrated trust and reputation model for open multi-agent systems. *Autonomous Agents and Multi-Agent Systems*, 13(2):119–154, 2006.

5. N. R. Jennings. Commitments and conventions: The foundation of coordination in multi-agent systems. *The knowledge engineering review*, 8(03):223–250, 1993.
6. A. Josang, R. Ismail, and C. Boyd. A survey of trust and reputation systems for online service provision. *Decision Support Systems*, 43(2):618–644, 2007.
7. J. E. Kittock. Emergent conventions and the structure of multi-agent systems. In *Proc. of Santa Fe Institute Complex Systems Summer School VI*, pages 1–14, 1993.
8. J. Kleinberg. Navigation in a small world. *Nature*, 406(3):845, September 2000.
9. M. W. Macy and J. Skvoretz. The evolution of trust and cooperation between strangers. *American Sociological Review*, 63:638–660, 1998.
10. S. Mahmoud, N. Griffiths, J. Keppens, and M. Luck. Efficient norm emergence through experiential dynamic punishment. In *Proc. of the 20th European Conference on Artificial Intelligence*, pages 576–581, 2012.
11. J. Morales, M. López-Sánchez, and M. Esteva. Using Experience to Generate New Regulations. In *Proc. of the 22th Int. Joint Conf. on Artificial Intelligence*, pages 307–312, 2011.
12. M. Nowak and K. Sigmund. The dynamics of indirect reciprocity. *Journal of theoretical biology*, 194(4):561–74, October 1998.
13. M. A. Nowak and K. Sigmund. Evolution of indirect reciprocity. *Nature*, 437(7063):1291–1298, 2005.
14. A. A. Pirzada and C. McDonald. Trust Establishment In Pure Ad-hoc Networks. *Wireless Personal Communications*, 37(1-2):139–168, 2006.
15. S. D. Ramchurn, D. Huynh, and N. R. Jennings. Trust in multi-agent systems. *The Knowledge Engineering Review*, 19(01):1–25, April 2005.
16. J. Sabater, M. Paolucci, and R. Conte. Reputa: Reputation and image among limited autonomous partners. *Journal of Artificial Societies and Social Simulation*, 9(2):3, 2006.
17. N. Salazar, J. A. Rodríguez-Aguilar, and J. L. Arcos. Robust coordination in large convention spaces. *AI Communications*, 23(4):357–372, 2010.
18. S. Sen and S. Airiau. Emergence of norms through social learning. In *Proc. of the 20th Int. Joint Conf. on Artificial Intelligence*, pages 1507–1512, 2007.
19. Y. Shoham and M. Tennenholtz. On the emergence of social conventions: modeling, analysis, and simulations. *Artificial Intelligence*, 94(1-2):139–166, July 1997.
20. D. Villatoro, S. Sen, and J. Sabater-Mir. Topology and memory effect on convention emergence. In *Proc. of the 2009 Int. Conf. on Web Intelligence and Intelligent Agent Technologies*, pages 233–240, 2009.
21. A. Walker and M. Wooldridge. Understanding the emergence of conventions in multi-agent systems. In *Proc. of the 1st Int. Conf. on Multi-Agent Systems*, pages 384–389, 1995.
22. H. P. Young. The economics of convention. *The Journal of Economic Perspectives*, 10(2):105–122, 1996.