

# Towards Modeling Trust Based Decisions: A Game Theoretic Approach

S.Vidyaraman, M.Chandrasekaran and S.Upadhyaya  
{vs28,mc79,shambhu}@cse.buffalo.edu

University at Buffalo, Buffalo NY 14260, USA  
Phone: 716-645-3180, Fax: 716-645-3464

**Abstract.** Current trust models enable decision support at an implicit level by means of thresholds or constraint satisfiability. Decision support is mostly included only for a single binary action, and does not explicitly consider the purpose of a transaction. In this paper, we present a game theoretic model that is specifically tuned for decision support on a whole host of actions, based on specified thresholds of risk. As opposed to traditional representations on the real number line between 0 and +1, Trust in our model is represented as an index into a set of actions ordered according to the agent's preference. A base scenario of *zero trust* is defined by the equilibrium point of a game described in normal form with a certain payoff structure. We then present the *blind trust* model, where an entity attempts to initiate a trust relationship with another entity for a one-time transaction, without any prior knowledge or recommendations. We extend this to the *incentive trust* model where entities can offer incentives to be trusted in a multi-period transaction. For a specified risk threshold, both models are analyzed by using the base scenario of *zero trust* as a reference. Lastly, we present some issues involved in the translation of our models to practical scenarios, and suggest a rich set of extensions of the generalized game theoretic approach to model decision support for existing trust frameworks.

**Keywords:** Decision Support, Game Theory, Incentives, Risk, Trust

## 1 Introduction

The three dominant characteristics of trust are *vulnerability*, *risk* and *expectation* (or *uncertainty*). All trust models encompass these characteristics and present definitions, representations, evaluations and operations on the notion of trust (see [1–5] and the references therein). Decision support for trust models and frameworks must involve an accurate estimation of the uncertainty of other agent's actions. The level to which an agent is willing to tolerate the loss due to the uncertainty is the risk threshold. Most trust models/frameworks enable decision support based on threshold values or constraint satisfiability (e.g., *automated trust negotiations* first initiated by Winsborough et al. [6] and later extended in [7–9]) or some aggregation metric of recommendations (e.g., Fuzzy metrics [10]) based on past history, like recommender systems (RS) [11]. In most

of the previous works that present generic trust models, the decision making criteria, i.e., the translation from trust to action, is left to the agent, and rightly so, because such translations are usually context dependent. Decision support is not explicitly embedded into the trust model; rather the agent is expected to make decisions for a single action based on thresholds or constraints, depending on the model. Such threshold or constraint specification is for a single binary action and is not applicable when the agent has a multitude of action choices. The traditional trust values of 0 to +1 are not particularly conducive towards the direct translation of trust to a multitude of actions; an additional mapping function is required for decision support.

In this paper, we present a model of trust based decisions using a game theoretic approach. In our model, agents have a multitude of action choices and interact with other agents with some (possibly zero) trust. The trust of an agent (on other agents) is represented as an index into the action set of the agent. Thus, the very value of trust enables decision support, even at the level of the abstract model. In other words, the trust value describes the action to be initiated in an interaction. We extend the work of other trust models by going beyond defining trust notions; taking our intuition from *automated trust negotiations* [6], our model assumes that every trust interaction has a purpose, and thus, *both* (and in general *all*) interacting agents must have something to gain at the end of an interaction. With this purpose in mind, we first present a base scenario/game where interacting agents do not trust each other, and thus play at their equilibrium point in game theoretic terms. We then present the basic *blind trust* model, which is a one-time transaction between two agents. Here, an agent is assumed to trust another agent for reasons outside the scope of the model; no assumptions on the application domain are made, neither are reasons for trusting an agent provided. At this point, we define the notion of a *desired action* (in game theoretic terms) that formalizes the expectation of a trusting agent. We then define two types of risk in the *blind trust* model and evaluate the number of rounds of sequential game play a trusting agent (or truster) may expect to play for a given risk tolerance. Next we consider the purpose of a trust interaction and present the *incentive trust* model, where agents can provide an incentive to other agents in order to adhere to a minimal trust level that is established or agreed upon in advance. The *incentive trust* model also provides an advantage to a newly entering agent in the transaction who has no history; instead of starting from a default low trust value, the agent may quickly build up its reputation by offering incentives. In both these models, we use the base scenario as a reference point for determining the loss of an agent when trust is misplaced or violated. This loss is used to derive metrics for estimating the risk faced by an agent. Finally, we present a rich set of models and frameworks to which the general game theoretic approach may be extended. To the best of our knowledge, such explicit decision support and analysis in trust models through a game theoretic approach has not been done so far.

In this paper, we do not explicitly define the notions of trust (and distrust), (atomic) purpose of the trust relationship, etc. These notions have been well

defined in previous works [4, 1, 5] (also see [3, 2] and the references therein for a listing of previous works and models); indeed, the core notion of trust could stem from any of the prior works that not only define fundamental concepts well, but also provide means for evaluating and performing other operations (like comparison) on trust through the history of past transactions. Unlike previous approaches, the purpose of this work is to provide a game theoretic model of trust based decisions. This work represents a natural progression of existing trust models to provide explicit decision support for agents with a multitude of action choices. By proper mechanism design [12, 13], game theoretic models can also subsume other models to provide a well analyzed decision support theory.

### 1.1 Summary of Contributions

The contributions of this paper are in the theoretical realm of trust and are summarized as follows.

1. We present a game theoretic model for enabling trust based decision support by defining trust as an index into the action set of a trusting agent.
2. A *blind trust* model is presented, where agents engage in repeated games; two types of risk are defined and the number of rounds an agent may expect to play is analyzed for a given risk threshold.
3. An *incentive trust* model, where all interacting agents stand to gain at the end of the transaction is presented.
  - (a) Sufficient conditions are derived for both the agent offering the incentive and the trusting agent/truster in order for the interaction to be successful.
  - (b) This model offers a mechanism for a new agent to start an interaction with a high level of trust instead of the default low value by offering incentives.
4. We present directions for the translation of the game theoretic models for practical applications and suggest potentially rich areas of future works.

The rest of this paper is organized as follows. Section 2 describes the related work; Section 3 presents a brief background and intuitive description of the game theoretic approach of our model. Section 4 presents the base scenario of *zero trust*, the *blind trust* and *incentive trust* models. Section 5 presents some of the issues involved in applying the game theoretic models to practical applications. Concluding remarks and directions for future works are presented in Section 6.

## 2 Related Work

The work in this paper primarily focuses on enabling decision support for agents operating on the notion of trust. By its very definition [1–5], trust implies a certain amount of risk due to uncertainty in the interacting agents decision criteria. Trust and risk have both been used to decide or optimize the effective payoff [14] or lower the expected loss [15]. Their relationship towards decision

support has been investigated in [16]. An important factor of reciprocity in terms of trust has been experimentally investigated in [17].

Game theoretic descriptions and analysis of trust have been investigated by [18]; some games that cannot be represented in a normal form have been investigated experimentally [19]. Decision support based on trust has been investigated for electronic transactions [20], trust negotiations [7, 8, 21, 9, 6], etc. In fact, almost all trust models incorporate some decision criteria at least implicitly; however most of them are binary or threshold based, in that an action may be initiated if a certain constraint is satisfied or the trust value is above a certain threshold. In this work, we explicitly consider the purpose of any trust transaction and assume that all interacting entities gain at the end of the transaction. Our work is similar to game theoretic modeling of an auction marketplace, where agents choose actions with optimal payoffs. The work by Lam et al. [14] discusses trade in open marketplaces using trust and risk, and is the closest to the work in this paper. Our work substantially builds on the work of previous trust models by providing a model for trust based decision support, particularly in situations where an agent has a multitude of actions to choose from; agents can not only decide which action to take based on their trust value, but also evaluate the number of rounds of interaction (game) they can engage for a given threshold of risk.

### 3 Background and Overview

In order to make this paper self sufficient, we first give a brief descriptive background of a 2-player game and its corresponding Nash Equilibrium. We then present an overview of our notion of trust in a 2-player game. A 2-player play game is defined as a game between the players denoted as  $P_1$  and  $P_2$ . Each player is required to operate simultaneously over an action space. On the completion of an action, both players receive a payoff or a reward depending on the actions chosen by themselves and the other player. Games where the players have opposing goals are called non-cooperative games. The goal of each player is to maximize his or her own payoff. Towards this goal, each player develops a strategy over his/her actions spaces, thereby ensuring the best payoff in the game. Assuming that both players know each others action space and their corresponding payoffs, their strategy will be to choose the *best response* action for the other player's best strategy. Furthermore, assuming both players are perfectly rational, i.e., they both can efficiently compute each others best strategy recursively, their final action will be one from which each can hope to gain nothing by deviating unilaterally. Intuitively, such a 'final' action results in 'equilibrium' for both players; in order to maximize their payoffs, each player only need to play that particular action, regardless of how many times the game is played. Such an action profile is called a Nash Equilibrium.

In the interest of keeping the overview concise and intuitive, a number of details have been omitted; e.g., not all games have a single action strategy that results in equilibrium; there is usually a probability distribution on the action

space that also results in equilibrium (such an action profile is called a mixed strategy, while the previous single action equilibrium is called a pure strategy). The reader is referred to [22, 23] for a more detailed exposition on game theory.

We now describe the incorporation of trust into a standard 2-player non-zero sum game. In any game in the equilibrium situation, players have no incentive to deviate from their chosen actions. We formulate our model as follows. Consider the players ( $P_1$ ,  $P_2$ ) whose equilibrium actions are  $(a_1, b_1)$  with a payoff of (50, 50). If  $P_1$  were to choose some other action  $a_2$ , then in a typical game, there exists a *best response* action  $b_2$  for player  $P_2$  such that the effective payoff is (40,100), i.e.,  $P_2$ 's payoff would be larger than the equilibrium payoff and  $P_1$ 's payoff would be smaller (maybe negative too, but we assume positive payoffs). Now assume that there exists an action  $b_3$ , called *the desired response*, for player  $P_2$  such that the payoff is (60, 70), i.e., both players stand to gain from the equilibrium payoff, but player  $P_2$  stands to gain lesser than the optimum *best response* for action  $a_2$  (which happens to be action  $b_2$ ). Thus, player  $P_1$  is said to trust player  $P_2$ , if on playing  $a_2$ , there is an expectation that  $P_2$  would respond with  $b_3$  instead of  $b_2$ , thereby leading to a profit on both sides, but *not necessarily the maximum allowable* for player  $P_2$ . Now imagine a continuum of such actions  $a_k, a_{k+1}$ , etc., for  $P_1$  such that  $P_2$  can respond with actions  $b_k, b_{k+1}$ , etc., such that their payoffs are increasingly better than the equilibrium payoff, but  $P_2$ 's payoff is lesser than the best response actions  $b_{best-response}$  to  $a_k, a_{k+1}$ , etc. Then  $P_1$ 's trust in  $P_2$  is the index  $k$  into his action profile. If the action profiles are suitably ordered, an increasing index value indicates an increasing level of trust.

The intuition behind our model is simple: an act of trust implies, amongst other things, (a) a potential *vulnerability* on the part of the trusting agent, (b) a threshold of *risk* the trusting agent is willing to tolerate and (c) an uncertainty (or *expectation*) on the response of the other agent: i.e., the three dominant characteristics, *vulnerability*, *risk* and *expectation/uncertainty* have to be embedded into the model. The *vulnerability* on the part of the trusting agent is expressed by its deviation from the equilibrium play. The extent to which the trusting agent is willing to expose itself to the vulnerability is the *risk*. The responding agent may initiate the *best response* (and hence violate the trust placed in him) or initiate the *desired response*; thus, there is an *uncertainty* on the response type. The equilibrium point/action profile forms the *Base Scenario*, where the players do not trust each other.

### 3.1 *Blind Trust and Incentive Trust models*

The *blind trust* model presents the scenario where an agent trusts another agent with no assurances or guarantees. In this model, the vulnerability faced by an agent when trusting another agent is expressed in terms of the possible loss of payoff for one round. Then, given the risk an agent is willing to take (the maximum vulnerability), we analyze the number of rounds of the game the player may expect to play before getting back to the equilibrium play. Thus far, the trust of a single player is unconditional; we now introduce the *incentive trust* model

where users can trade a predefined amount of their payoffs before the initiation of single round of the game. We now consider the interactions between players and their decision criteria when there is an expectation of a return at a later point in time, i.e., there is a *purpose* to the entire transaction and all involved entities expect to gain something at the end of the transaction. Many trust relationships fall under this category. Consider a customer who interacts with a service provider; the customer may pay a premium for a service that he expects at a later time. Thus, the customer can be said to trust the service provider to keep his end of the contract. Such a scenario is also applicable in the security field, in *Automated Trust Negotiation* mechanisms and protocols. Trust Negotiation was first introduced by Winsborough [6] as a means for agents to negotiate their trust with other agents in a heterogeneous environment. A commonly quoted example is the interaction between a potential customer Alice and a service provider Bob. The simplified scenario is as follows: a customer Alice wishes to make a purchase from Bob, a service provider, but is initially unwilling to provide any means of authentication (Drivers License, Credit Card Number, etc.). Bob provides Alice with a certificate from the Better Business Bureau (BBB) stating that he is indeed a service provider with a certain standing. The BBB certificate is verified by Alice and it ‘helps’ her to make a trust decision about Bob; she provides her Resellers License/Credit Card number to Bob to make a purchase. Bob verifies her license/CC number and completes the transaction. This example is illustrative of the general trust transaction: (a) there is a purpose to the transaction and (b) both (and in general all) agents stand to gain at the end of the transaction. We analyze these scenarios and present decision making criteria for specified risk thresholds.

### 3.2 A Note on Mechanism Design

Finally, we conclude with a note on the concept of mechanism design [12, 13] in game theory. Loosely speaking, mechanism design can be viewed as a technique for designing a game so that *rational* and *selfish* agents do not have an incentive to deviate from the desired behavior of the game designer. Proper mechanism design maps the desired behavior of agents to the equilibrium play so that no agent can gain by deviating from the equilibrium point. From a purely game theoretic viewpoint, assuming agents are selfish and rational, the equilibrium play is desirable. However, in our model, we stipulate a deviation from the equilibrium play, towards a play that leads to a greater, but not necessarily maximum payoff, in order to embed the notion of *vulnerability* of an agent. Thus, there is no stipulation for an agent to adhere to a specific action set; if there were one, we would not be able to embed the notion of *vulnerability* and *uncertainty*; indeed, a stipulation of any kind would imply determinism, which runs contrary to free will or choice of an agent. However, there are game theoretic constructs like satisficing game theory [24, 25], amongst others, that can model users (as opposed to automated agents). We mention such games in Section 6 on extensions of this model; they are not considered in this paper. Herein, we shall use the

terms agent, user or player interchangeably; they imply the same entity unless specifically mentioned otherwise.

## 4 Trust Based Decision (TBD) Model

### 4.1 Base Scenario: *Zero Trust*

We now fix the base scenario, which is a 2-player game with a Nash Equilibrium. Our notations closely follow standard game theoretic expositions as in [26]; although our models consider only 2 players, the notations are kept generic. The trust game  $G_\tau$  is defined by:

$$G_\tau = \{N, (A_i)_{i \in N}, (u_i)_{i \in N}\} \quad (1)$$

where  $N$  is the set of players,  $A_i$  is the action space of player  $i$  and  $u_i$  is the payoff function for player  $i$ , defined as  $u_i: \mathbf{A} \rightarrow \mathbb{R}$ , where  $\mathbf{A} = \times_{i \in N} A_i$  and  $\mathbb{R}$  is the set of real numbers. We assume that  $N$  and  $A_i$  are finite sets. Player  $i$  has at his disposal the actions  $a_i \in A_i$ . We denote the action spaces of all other players other than  $i$  as  $A_{-i} = \times_{j \in N \setminus \{i\}} A_j$  and a single element as  $a_{-i} \in A_{-i}$ . The repeated rounds of the game  $G_\tau$  are referred to as the *supergame*, which consists of a finite sequence of the game  $G_\tau$ , where the players choose the actions  $a_i(t) \in A_i$  at time instant  $t$ . For a sequence of  $k$  plays, we denote the history of player  $i$  by  $H_i(k) = \{a_i(1), a_i(2), \dots, a_i(k)\} \forall a_i(\cdot) \in A_i$  and each element of  $H_i(k)$  by  $h_i(k)$ . The payoff of player  $i$  at the end of any round is given by  $u_i(a_i, a_{-i})$ . The best response action of player  $i$  is defined as  $b_i(a_{-i}) = \{a_i \in A_i : u_i(a_i, a_{-i}) \geq u_i(a_i^*, a_{-i}) \forall a_i^* \in A_i\}$ , i.e., given the plays of all the opponents  $a_{-i}$ , the action  $b_i(a_{-i})$  ensures the best payoff for player  $i$ . Hereafter,  $b_i(a_{-i})$  is denoted simply as  $b_i$ . We assume that the game's payoff structure allows for (at least) a single equilibrium point, at which the action profile of the players is  $(b_i, b_{-i})$ . Thus the single round payoff of the player  $i$  is given by  $u_i(b_i, b_{-i})$ . Intuitively, the cumulative payoff of a sequence of  $k$  plays is  $k \cdot u_i(b_i, b_{-i})$ . However, for the game  $G_\tau$ , similar to [26], we define the cumulative payoff to be a discounted one, where the weights of the payoffs of older sequences are progressively lesser.

**Definition 1.** *The discounted cumulative payoff in the game  $G_\tau$ , of player  $i$  over a sequence of  $k$  play's, discounted by a factor of  $\delta \in (0, 1)$ , is defined as:*

$$C_i(\delta, k) = (1 - \delta) \sum_{m=1}^k \delta^{k-m} u_i(a_i(m), a_{-i}(m)) \quad (2)$$

Unless otherwise specified, we shall denote  $C_i(\delta, k)$  as  $C_i(\delta)$ . This formulation places greater relevance to the most recent play (the  $k^{th}$  play) and progressively decreases the payoff of the past plays. From a trust game viewpoint, this is intuitive; the closer  $\delta$  is to 1, the greater the relevance to the most recent play (due to the factor  $\delta^{k-m}$ ). Note that  $\delta \in (0, 1)$  and does not ever assume the

value of 0 or 1. However, if we set  $\delta = 0$ ,  $C_i(\delta) = 1$ ; this can be interpreted as setting no relevance at all to any of the plays, and hence the payoff incurred at any stage is a constant: in the context of the model, setting no relevance to the plays makes no sense and hence such a setting is invalid.<sup>1</sup> In the base scenario, the action taken by the players are the best responses to the action spaces of other players; hence the cumulative payoff of player  $i$  is obtained when  $a_i(\cdot) = b_i$  and  $a_{-i}(\cdot) = b_{-i}$  in Eq. 2. We denote this best response cumulative payoff as  $C_i^*(\delta)$ . The two values of per round payoff,  $u_i(b_i, b_{-i})$  and cumulative discounted payoff  $C_i^*(\delta)$  are used to refer to the *Base Scenario* with no trust.

## 4.2 Blind Trust

In this model, player  $i$  wants to initiate a trust relationship with the remaining players. His goal is now to obtain the *desired response* from the remaining players as opposed to the *best response*  $b_{-i}$  played by the remaining players in the *Base Scenario*.

**Definition 2.** For an action  $a_i \in A_i$ , the desired response  $d_i(a_i)$  of player  $i$  is defined as an action from the set  $A_{-i}$  that increases the payoff of all players from the equilibrium payoff, but does not provide the maximum possible payoff to all the players other than possibly player  $i$ .

$$\begin{aligned} d_i(a_i) = \{a_{-i} \in A_{-i} : u_i(a_i, a_{-i}) \geq u_i(b_i, b_{-i}), \\ u_{-i}(a_i, b_{-i}) \geq u_{-i}(a_i, a_{-i}) \geq u_i(b_i, b_{-i})\} \end{aligned} \quad (3)$$

The desired response of player  $i$  is also denoted simply as  $d_i$ , where the action  $a_i$  is understood to have been initiated. Note that  $d_i(a_i) \in A_{-i}$  and is not necessarily unique. Depending on the application domain, there may exist multiple  $d_i(a_i)$ ; however, their existence does not affect our model from a decision theoretic viewpoint.

**Definition 3.** The index value  $\tau$  of actions  $a_i^\tau \in A_i$  in a strictly increasing ordering given by  $u_i(a_i^1, d_i) \leq u_i(a_i^2, d_i) \leq \dots \leq u_i(a_i^\tau, d_i) \leq \dots \leq u_i(a_i^T, d_i)$  is defined as the trust that player  $i$  places on the remaining players.

Note that the desired action  $d_i$  for  $a_i^\tau$  is *not* (necessarily) the same for all index values. In situations where the context is clear, the symbol  $\tau$  is used to represent the trust of player  $i$ ; in more generic terms,  $\tau(i \rightarrow -i)$  represents the cumulative trust of player  $i$  on the remaining players, while  $\tau(i \rightarrow j)$  represents the trust of player  $i$  on player  $j$ . In the interest of keeping the formulation generic and extensible, the notations of  $i$  and  $-i$  have been used; herein, we shall restrict ourselves to  $N = 2$ , i.e., there are two players in the game  $G_\tau$ ; thus  $i \in \{1, 2\}$  ( $-i$  denotes the ‘other’ player). In the formulations that follow, replacing  $i = 1$  and  $-i = 2$  represents the model for the two player scenario. The trust value  $\tau \in \{1, T\}$ , where  $T$  is the maximum index. As a technical device, we may also

<sup>1</sup> We thank an anonymous reviewer for bringing out this point.



include a zero value in  $\tau$  where the index value of zero is associated with the best response action (and hence no trust).

In the basic *blind trust* scenario, player  $i$  assigns a trust value  $\tau = 1$ , and waits for the response from player  $-i$ . From the game theoretic viewpoint, this game is a turn based game, where player  $-i$  knows the action taken by player  $i$  before his turn to play. We call this model a *blind trust* model since there is no prior communication between the two players for a contractual agreement on the action set, etc. Player  $i$  blindly trusts player  $-i$  and hopes for a reciprocation. At this point, player  $-i$  may reciprocate by initiating the desired response  $d_i$  or the best response  $b_{-i}$ . The initiation of the desired response indicates the beginning of a trust relationship.

From this basic *blind trust* Model, we wish to address several questions. First, given that an agent (player  $i$ ) wishes to initiate a trust relationship, what is the vulnerability faced by player  $i$ ? Secondly, assume that player  $i$  initiates a *blind trust* relationship in the hopes of a future collaboration. Initially, player  $-i$  may simply act in a ‘rational’<sup>2</sup> manner and initiate the best response action, but may eventually reconsider or ‘understand’ that player  $i$  wishes to initiate a trust relationship. The reasons for the establishment and evolution of the trust relationship are contextual and application/domain dependent, and are hence outside the scope of this paper. However, the relevant question is: given the amount of risk (maximum vulnerability) that player  $i$  is willing to tolerate, what is the number of rounds of play that player  $i$  may expect to play? Towards answering this question, we first define two types of risk and then evaluate the expected number of rounds.

**Instantaneous Per-Round and Cumulative Risk:** The risk faced by a player  $i$  are categorized into two types: the *instantaneous per-round risk* and the *cumulative  $k$ -stage risk*.

**Definition 4.** *The instantaneous per-round risk  $\rho_i$  of player  $i$  when initiating action  $a_i^\tau$  with a trust  $\tau$  on player  $-i$  is defined as the ratio of the difference between the equilibrium payoff  $u_i(b_i, b_{-i})$  and the best response  $u_i(a_i^\tau, b_{-i})$  to the equilibrium payoff.*

$$\rho_i(\tau) = \left(1 - \frac{u_i(a_i^\tau, b_{-i})}{u_i(b_i, b_{-i})}\right) \quad (4)$$

Note that, by the very definition of best response actions,  $u_i(a_i^\tau, b_{-i}) \leq u_i(b_i, b_{-i})$  (otherwise,  $a_i^\tau = b_i$ ). Thus,  $\rho_i(\tau)$  is the risk faced by the player  $i$  when trusting player  $-i$  with a value of  $\tau$ . Intuitively, the *simplified cumulative  $k$ -stage risk* is  $k \cdot \rho_i(\tau)$ , assuming that the player  $-i$  plays the best response for all the  $k$  sequences of the game.

Recall that we have defined a cumulative discounted payoff for the game  $G_\tau$  in definition 1, Eq. 2, i.e., the payoff of the  $k^{th}$  round is discounted by a factor  $(1 - \delta)$ . Towards this, a discounted cumulative  $k$ -stage risk is defined, similar to definition 4. We first derive the discounted cumulative  $k$ -stage equilibrium

---

<sup>2</sup> ‘Rational’ in the game theoretic sense, not in the context of the application domain.

payoff and best response payoff. The *discounted cumulative k-stage equilibrium payoff* of player  $i$  can be derived by substituting  $a_i(m) = b_i$  and  $a_{-i}(m) = b_{-i} \forall m = \{1, 2, \dots, k\}$  in Eq. 2.

$$C_i^{eq}(\delta) = (1 - \delta) \sum_{m=1}^k \delta^{k-m} u_i(b_i, b_{-i}) = u_i(b_i, b_{-i})(1 - \delta^k) \quad (5)$$

Note that the discounted equilibrium payoff is in fact the same as the equilibrium payoff for any single round if  $\delta \rightarrow 0$ , in which case, it is almost *independent* of the number of rounds over which the game is played. Similarly, the *discounted cumulative k-stage best response payoff* of player  $i$  can be derived by substituting  $a_i(m) = a_i^\tau$  and  $a_{-i}(m) = b_{-i} \forall m = \{1, 2, \dots, k\}$  in Eq. 2.

$$C_i^{best}(\delta) = (1 - \delta) \sum_{m=1}^k \delta^{k-m} u_i(a_i^\tau, b_{-i}) = u_i(a_i^\tau, b_{-i})(1 - \delta^k) \quad (6)$$

**Definition 5.** The *discounted cumulative k-stage risk*  $\sigma_i(\tau, k)$  of player  $i$  over a sequence of  $k$  plays of the game  $G_\tau$  for the histories  $H_i(k)$  and  $H_{-i}(k)$  is defined as the ratio of the difference between the equilibrium payoff  $C_i^{eq}(\delta)$  and the best response payoff  $C_i^{best}(\delta)$  to the equilibrium payoff.

$$\sigma_i(\tau, k) = \frac{C_i^{eq}(\delta) - C_i^{best}(\delta)}{C_i^{eq}(\delta)} = \rho_i(\tau) \quad (7)$$

Usually, the cumulative risk is evaluated until player  $-i$  plays the desired action  $d_i(a_i^\tau)$ .

**Expected Number of Plays:** Consider the situation when player  $i$  wishes to initiate a trust relationship and initiates the action  $a_i^\tau$  instead of  $b_i$ . Assume that the maximum amount of risk the player  $i$  is willing to undertake is  $r_i$ . We now evaluate the number of rounds that player  $i$  may expect to play.

Consider the *simplified cumulative k-stage risk*,  $k, \rho_i(\tau)$ . In this case, the maximum number of rounds player  $i$  can afford to play is  $k_{max} = \frac{\rho_i(\tau)}{r_i}$ . Now assume that at any round of the  $k$  stages, the probability that player  $-i$  switches from  $b_{-i}$  to  $d_i$  is  $p$ . This probability is available to the player  $i$  through some context specific mechanism, also called a *belief* in game theoretic literature. The probability that player  $-i$  switches from  $b_{-i}$  to  $d_i$  at the  $k^{th}$  round is  $p(1-p)^{k-1}$ . Thus the number of rounds player  $i$  may expect to play is  $\sum_{m=1}^{k_{max}} mp(1-p)^{m-1}$ . Thus, considering the *simplified cumulative k-stage risk*, the number of rounds player  $i$  may expect to play is:

$$E[k] = \frac{1 - (k_{max} + 1)(1-p)^{k_{max}} + k_{max}(1-p)^{k_{max}+1}}{p} \quad (8)$$

where  $k_{max} = \frac{\rho_i(\tau)}{r_i}$ . Now, let's consider the *discounted cumulative k-stage risk*  $\sigma_i(\tau, k)$ , which is equal to  $\rho_i(\tau)$ . It can be observed trivially, that if the risks are discounted for past rounds of the game, then player  $i$  may continue to play the

game infinitely if  $r_i > \rho_i(\tau)$ , just one round if  $r_i = \rho_i(\tau)$  and may not play at all if  $r_i < \rho_i(\tau)$ ; i.e., the expected number of rounds is 1. Thus, for the discounted case, the specification of the player's *absolute risk* is not useful to determine the maximum number of rounds. Instead, we define the *discounted loss* (in payoff) for the player  $i$  for  $k$  rounds as:

$$L_i(k) = C_i^{eq}(\delta) - C_i^{best}(\delta) = (u_i(b_i, b_{-i}) - u_i(a_i^\tau, b_{-i}))(1 - \delta^k) \quad (9)$$

Let the loss that player  $i$  is willing to sustain in the trust initiative be  $l_i$ . Thus, the maximum number of rounds is given by:

$$L_i(k_{max}) = l_i \Rightarrow k_{max} = \frac{\log(1 - \frac{l_i}{u_i(b_i, b_{-i}) - u_i(a_i^\tau, b_{-i})})}{\log \delta} \quad (10)$$

Note that the loss  $l_i$  is not the absolute loss, but is the (maximum) an agent is willing to sustain relative to the best response action  $b_{-i}$ ; hence  $l_i < (u_i(b_i, b_{-i}) - u_i(a_i^\tau, b_{-i}))$ . Note also that  $k_{max}$  is inversely proportional to  $\log \delta$ , i.e., a player's maximum number of rounds depends on the extent to which he is willing to discount past payoffs (or equivalently, losses). The expected number of rounds, assuming that the probability that player  $-i$  switches from  $b_{-i}$  to  $d_i$  is  $p$ , is given by:

$$E[k] = \frac{1 - (k_{max} + 1)(1 - p)^{k_{max}} + k_{max}(1 - p)^{k_{max} + 1}}{p} \quad (11)$$

where  $k_{max}$  is given by Eq. 10.

**Compensatory Update Strategy:** Once the player  $-i$  responds with  $d_i$ , player  $i$  may update the value of  $\tau$ . Given that we are dealing with payoff values, we may use an update strategy that assigns  $\tau$  based on the payoff values. Our previous work [27] describes a *Compensatory Trust Model* (CTM), where the trust value may be updated as part of a compensation given to player  $-i$  based on his forgone payoff. We briefly describe the intuition behind the update strategy based on the CTM. When player  $-i$  initiates the desired action  $d_i$  at a particular round, denote his loss of payoff as  $l_{-i} = (u_{-i}(a_i^\tau, b_{-i}) - u_{-i}(a_i^\tau, d_i))$ , and the gain in payoff of player  $i$  as  $g_i = (u_i(a_i^\tau, d_i) - u_i(a_i^\tau, b_{-i}))$ . To 'share' the loss and gain equally, player  $i$  would have to transfer a payoff of  $l_{-i} + \frac{1}{2}(g_i - l_{-i})$  to player  $-i$ . This transfer may be made figuratively by updating the trust value proportionally, for some  $\delta \in (0, 1)$ ,  $\tau = \delta(l_{-i} + \frac{1}{2}(g_i - l_{-i}))$ , i.e., the trust update is proportional to the discount  $\delta$  of the payoffs. This is the 'compensation' that player  $i$  pays to player  $-i$ ; hence the name *Compensatory Update Strategy*. This scheme can be extended to include the risk faced by player  $i$  or the loss of payoff in a  $k$ -stage game. As mentioned before, we do not present new trust assignment and update methodologies; apart from the CTM update described above, any update mechanism may be used to assign and update trust.

### 4.3 Incentive Trust

The *incentive trust* model considers the purpose of a trust transaction and models the scenario where agents may provide incentives to be trusted. Consider

the game  $G_\tau$  with the players  $i$  and  $-i$ . The basic philosophy behind the trust transaction, unlike the *blind trust* model, is the expectation of a return in trust or value/payoff by a player. Towards this, we recast the game play described in [28] to fit into our model. The *incentive trust* model is a game played over three time periods. The play of the *incentive trust* model is described as follows:

**Initial Setup:** Player  $-i$  wants player  $i$  to (a) trust him with a level of  $\tau(i \rightarrow -i)$ , hereafter denoted as  $\tau$  and (b) initiate the action  $a_i^\tau$ . Towards this, player  $-i$  states  $p_{-i}$ , his *stated* probability that he will respond with the action  $d_i$ , the desired action instead of  $b_{-i}$ , the best response action.  $q_{-i}$  is the *true* probability that player  $-i$  will respond with the action  $d_i$ , the desired action instead of  $b_{-i}$ , the best response action. In a similar vein, let  $p_i$  be the *stated* probability (while  $q_i$  is the real probability) that player  $i$  will initiate action  $a_i^\tau$ . Both  $p_{-i}$  and  $p_i$  are public knowledge.

**Time Period 1:** Player  $-i$  transfers a payoff of  $f(p_{-i})$  as a support/proof of his commitment to respond with the desired action  $d_i$ , where  $f(\cdot)$  is a function whose specification is to be determined.

**Time Period 2:** Player  $i$ , on receiving the support of  $f(p_{-i})$  initiates the action  $a_i^\tau$ .

**Time Period 3:** Player  $-i$ , in turn, initiates the desired action  $d_i$ .

The expected payoff of player  $i$ ,  $w(p_{-i})$  can be split into the following components:

1. Player  $i$  gets a payoff of  $f(p_{-i})$  in time period 1.
2. Player  $i$  gets a payoff of *only*  $(1 - p_{-i})q_i u_i(a_i^\tau, b_{-i})$  if player  $-i$  cheats and responds with the best response  $b_{-i}$  instead of the desired response.
3. Player  $i$  gets a payoff of  $q_i p_{-i} u_i(a_i^\tau, d_i)$  if he acts with action  $a_i^\tau$  and player  $-i$  responds with the desired action.

Thus, the expected payoff of player  $i$  is given as:

$$w(p_{-i}) = f(p_{-i}) + q_i \{ (1 - p_{-i}) u_i(a_i^\tau, b_{-i}) + p_{-i} u_i(a_i^\tau, d_i) \} \quad (12)$$

Similarly, the expected payoff of player  $-i$ ,  $z(p_i)$  can be split into the following components:

1. If Player  $i$  initiates an action of  $a_i^\tau$ ,
  - (a) Player  $-i$  gets a payoff of  $p_i(1 - q_{-i})u_{-i}(a_i^\tau, b_{-i})$  if he responds with the best action  $b_{-i}$ .
  - (b) Player  $-i$  gets a payoff of  $p_i q_{-i} u_{-i}(a_i^\tau, d_i)$  if he responds with the desired action  $d_i$ .
2. If Player  $i$  does not initiate an action of  $a_i^\tau$ , then player  $-i$  loses a payoff of  $(1 - p_i)f(p_{-i})$ .

Thus, the expected payoff of player  $-i$  is given as:

$$z(p_i) = p_i(1 - q_{-i})u_{-i}(a_i^\tau, b_{-i}) + p_i q_{-i} u_{-i}(a_i^\tau, d_i) - (1 - p_i)f(p_{-i}) \quad (13)$$

**Analysis:** Consider the entire transaction and the purpose of the model: from the viewpoint of player  $i$ , we would like to obtain an optimal guarantee  $f(p_{-i})$  from player  $-i$ . Player  $-i$ , on the other hand, would like to ensure he at least has a non-zero payoff from the entire transaction. Since Player  $i$  has the first position advantage, given  $p_{-i}$ , he can fix the value of  $f(p_{-i})$  that he receives from player  $-i$ . The optimal value  $f_{opt}(p_{-i})$  is thus given by the first-order condition  $w'(p_{-i}) = 0$ .

$$f_{opt}(p_{-i}) = p_{-i}q_i\{u_i(a_i^\tau, b_{-i}) - u_i(a_i^\tau, d_i)\} + c \quad (14)$$

where  $c$  is a constant such that  $c \geq 0$ . If  $f(p_{-i})$  satisfies Eq. 14, then player  $i$  can expect a payoff of  $w(p_{-i}) = q_i u_i(a_i^\tau, b_{-i}) + c$ , independent of  $p_{-i}$ . From the viewpoint of player  $-i$ , the condition to be satisfied is  $z(p_i) \geq 0$ . Thus, we derive the condition for  $q_{-i}$ , taking into account the value of  $f(p_{-i})$  given in Eq. 14.

$$q_{-i} \geq \frac{p_i u_{-i}(a_i^\tau, b_{-i}) + (1 - p_i)c}{p_i(u_{-i}(a_i^\tau, b_{-i}) - u_{-i}(a_i^\tau, d_i)) + (1 - p_i)p_{-i}(u_i(a_i^\tau, b_{-i}) - u_i(a_i^\tau, d_i))} \quad (15)$$

Note that  $q_{-i}$  is not only dependent on  $p_i$ , but also on  $p_{-i}$ , on account of its dependence on  $f(p_{-i})$ . Since player  $-i$  is the initiator in this model, it is his responsibility to ensure a high enough value of  $q_{-i}$  in order to ‘break even’ in the transaction.

## 5 From Theory to Practice: Some Issues

The translation of game theoretic models to practical application scenarios is not without its problems [29]. Although the *blind trust* and *incentive trust* models capture the purpose of a trust transaction at an intuitive level, there are some issues that should be addressed before its translation to practical scenarios. We illustrate some of these issues and suggest directions for practical usage.

1. We have assumed that each user’s payoff is transferable to the other user; e.g., in the *incentive trust* model player  $-i$  initially transfers an amount of  $f(p_{-i})$ . Games where players can transfer their payoffs to others can be modeled by a class of games called Transferable Utility (TU) games, which also permit coalitions among groups of players. In practical scenarios, transferable payoffs have to be defined either in terms of monetary values or resources/services/QoS-guarantees, etc., depending on the application domain.
2. We have not defined the semantics of the game which would indicate the relevance of the equilibrium (or no trust) plays. The semantics of the game must be defined for the particular application domain where the model is applied to.
3. The most difficult part is the cohesion of the semantics of the game with the *zero trust* interaction. In practical scenarios, *zero trust* usually implies a lack of interaction and thus, no game plays at all. Bridging the game

semantics so that game plays occur in all scenarios is a challenging task; in particular, the process of mechanism design must accurately map to the scenarios so that the evaluation of iterated plays is relevant to the situation under consideration. Of particular interest here is to model games which also specify distrust.

4. The specification or formulation of the utility functions is the first step towards enabling decision support based on trust. Depending on the application domain, utility functions may be formulated in functional form or, if feasible, can be specified in discrete form for each action in the action set. The works [30, 31] provide methodologies for the construction of utility functions depending on the history of an agent's utility realization.
5. One of the pitfalls of game theoretic analysis is the existence of recursive reasoning about the other agent's knowledge [32]. The manifestation of this pitfall is best illustrated with the *incentive trust* model. For example, player  $i$  decides  $f(p_{-i})$  based on the value of  $p_{-i}$ , the stated probability of player  $-i$  that he will play the desired action  $d_i$  instead of the best response  $b_{-i}$ . Although, there is a constant  $c \geq 0$  involved in Eq. 12 which gives player  $i$  a leeway in choosing  $f(p_{-i})$ , it may be argued that player  $-i$  can quote a low value of  $p_{-i}$  and hence, the value of  $c$  must be chosen by player  $i$ , depending, amongst other factors, on player  $i$ 's belief in the value of  $p_{-i}$ . Such an analysis path usually leads to recursive reasoning and as illustrated by experimental evidence [32–34], is not advisable in modeling what essentially is a subjective concept, beyond maybe two or three levels. Obviously, if  $p_{-i} < \frac{1}{2}$ , player  $i$  would not even consider entering into the transaction; neither would player  $-i$  enter the transaction if  $f(p_{-i})$  is too high.

## 6 Conclusion and Future Work

The purpose of most trust models and frameworks is to provide some form of decision support for an agent or user. Previous works in trust models have focused on defining the notion of trust and evaluating it based on parametric representations, recommendations, etc. They provide decision support at an implicit level, either through threshold based schemes or constraint satisfiability. This paper provides a game theoretic approach to model trust based decisions. The model takes a holistic view of a trust interaction by considering the ultimate purpose of the transaction which may be spread over multiple periods. The very definition of trust as an index into the agent's action set provides decision support. Research in this direction has a potential to provide quantitative decision support for trust frameworks with a multitude of actions choices. Future work on this model comprise of two distinct paths: *extensions* and *practical applications*. Trust decisions that are taken based on the recommendations of other players can be incorporated into the model by means of an information structure. Furthermore, the generic game  $G_\tau$  assumes the existence of at least one equilibrium point; the model can be improved by defining games with predefined payoff structures [35] that imply the existence of a pure equilibrium point. Lastly, trust negotiations

that are made by people instead of automated agents have properties different from automated agents. Satisficing game theory provides a mathematical basis for user preferences/negotiations [24], which also account for the interests of users in the welfare of other users.

## References

1. Levien, R.: Advogato Trust Metric. Phd thesis, UC Berkeley (2003)
2. Kemal, B., Bruno, C., Andrew, S.T.: How to incorporate revocation status information into the trust metrics for public-key certification. In: ACM Symposium on Applied Computing, TRECK Track, ACM Press (2005)
3. Josang, A., Ismail, R., Boyd, C.: A survey of trust and reputation systems for online service provision. *Decision Support Systems* **43**(2) (2007) 618–644
4. Ji, M., Orgun, M.A.: Trust management and trust theory revision. *IEEE Transactions on Systems, Man and Cybernetics, Part A* **36**(3) (2006) 451–460
5. Ray, I., Chakraborty, S.: A vector model of trust for developing trustworthy systems. In: Computer Security ESORICS 2004. (2004) 260–275
6. Winsborough, W., Seamons, K., Jones, V.: Automated trust negotiation. Technical report, North Carolina State University at Raleigh (2000)
7. Winsborough, W.H., Li, N.: Towards practical automated trust negotiation. In: Proceedings of the Third International Workshop on Policies for Distributed Systems and Networks (Policy 2002), IEEE Computer Society Press (jun 2002) 92–103
8. Yu, T., Winslett, M.: Policy migration for sensitive credentials in trust negotiation. In: WPES '03: Proceedings of the 2003 ACM workshop on Privacy in the electronic society, New York, NY, USA, ACM Press (2003) 9–20
9. Yu, T., Winslett, M., Seamons, K.E.: Interoperable strategies in automated trust negotiation. In: CCS '01: Proceedings of the 8th ACM conference on Computer and Communications Security, New York, NY, USA, ACM Press (2001) 146–155
10. Aringhieri, R., Damiani, E., Vimercati, S.D.C.D., Paraboschi, S., Samarati, P.: Fuzzy techniques for trust and reputation management in anonymous peer-to-peer systems. *Journal of the American Society for Information Science and Technology* **57**(4) (2006) 528–537
11. Avesani, P., Massa, P., Tiella, R.: A trust-enhanced recommender system application: Moleskiing. In: SAC '05: Proceedings of the 2005 ACM symposium on Applied computing, New York, NY, USA, ACM Press (2005) 1589–1593
12. Asselin, F., Jaumard, B., Nongailard, A.: A technique for large automated mechanism design problems. In: IAT '06: Proceedings of the IEEE/WIC/ACM international conference on Intelligent Agent Technology, Washington, DC, USA, IEEE Computer Society (2006) 467–473
13. Dash, R., Ramchurn, S., Jennings, N.: Trust-based mechanism design. In: Proceedings of the Third International Joint Conference on Autonomous Agents and MultiAgent Systems, ACM (2004) 748–755
14. Lam, Y.H., Zhang, Z., Ong, K.L.: Trading in open marketplace using trust and risk. In: IAT '05: Proceedings of the IEEE/WIC/ACM International Conference on Intelligent Agent Technology, Washington, DC, USA, IEEE Computer Society (2005) 471–474
15. Cvrcek, D., Moody, K.: Combining trust and risk to reduce the cost of attacks. In: Proceedings of the Third Annual Conference on Trust Management (iTrust 2005). Volume 3477 of LNCS., Springer-Verlag (May 2005) 372–383

16. Josang, A., Presti, S.: Analysing the relationship between risk and trust. In: Proceedings of the Second International Conference on Trust Management, Oxford (2004)
17. Cox, J.C.: How to identify trust and reciprocity. *Games and Economic Behavior* **46**(2) (2004) 260–281
18. Baras, J.S., Jiang, T.: Cooperative games, phase transitions on graphs and distributed trust in manet. In: 43rd IEEE Conference on Decision and Control. Volume 1. (2004) 93–98
19. Wu, D.J., Kimbrough, S.O., Fang, Z.: Artificial agents play the 'mad mex trust game': a computational approach. In: Proceedings of the 35th Annual Hawaii International Conference on System Sciences (HICSS), Hawaii (2002) 389–398
20. Josang, A.: Trust-based decision making for electronic transactions. In: Fourth Nordic Workshop on Secure IT Systems (NORDSEC'99), Stockholm, Sweden, Stockholm University Report 99-005. (1999)
21. Chen, W., Clarke, L., Kurose, J.F., Towsley, D.F.: Optimizing cost-sensitive trust-negotiation protocols. In: INFOCOM. (2005) 1431–1442
22. Davis, M.: *Game Theory: A nontechnical introduction*. Dover (1983)
23. Luce, R.D., Raiffa, H.: *Games and Decisions*. Dover (1989)
24. Archibald, J.K., Hill, J.C., Johnson, F.R., Stirling, W.C.: Satisficing negotiations. *IEEE Transactions on Systems, Man, and Cybernetics, Part C* **36**(1) (2006) 4–18
25. Candale, T., Sen, S.: Effect of referrals on convergence to satisficing distributions. In: AAMAS '05: Proceedings of the fourth international joint conference on Autonomous agents and multiagent systems, New York, NY, USA, ACM Press (2005) 347–354
26. Carmona, G.: A strong anti-folk theorem. *International Journal of Game Theory* **34**(1) (April 2006) 131–151
27. Vidyaraman, S., Upadhyaya, S.: A trust assignment model based on alternate actions payoff. In: Proceedings of the Fourth International Conference on Trust Management (iTrust '06), Volume 3986, Pisa, Italy (2006) 339–353
28. Huberman, B.A., Wu, F., Zhang, L.: Ensuring trust in one time exchanges: solving the qos problem. *NETNOMICS* **7**(1) (2005) 27–37
29. Mahajan, R., Rodrig, M., Wetherall, D., Zahorjan, J.: Experiences applying game theory to system design. In: PINS '04: Proceedings of the ACM SIGCOMM workshop on Practice and theory of incentives in networked systems, New York, NY, USA, ACM Press (2004) 183–190
30. Restificar, A., Haddawy, P.: Constructing utility models from observed negotiation actions. In: Proceedings of the Eighteenth Joint International Conference on Artificial Intelligence. (2003) 1404–1405
31. Afriat, S.: The construction of utility functions from expenditure data. Cowles Foundation Discussion Papers 177, Cowles Foundation, Yale University (1964)
32. Colman, A.M.: Depth of strategic reasoning in games. *Trends in Cognitive Sciences* **7**(1) (2003) 2–4
33. Hedden, T., Zhang, J.: What do you think i think you think?: Strategic reasoning in matrix games. *Cognition* **85**(1) (2002) 1–36
34. Stahl, D.O., Wilson, P.W.: Experimental evidence on players' models of other players. *Journal of Economic Behavior & Organization* **25**(3) (1994) 309–327
35. Monderer, D., Shapley, L.S.: Potential games. *Games and Economic Behavior* **14**(1) (1996) 124–143