

A Discriminatory Rewarding Mechanism for Sybil Detection with Applications to Tor

Asim Kumar Pal, Debabrata Nath and Sumit Chakraborty

Abstract—This paper presents an economic game for sybil detection in a distributed computing environment. Cost parameters reflecting impacts of different sybil attacks are introduced in the sybil detection game. The optimal strategies for this game in which both sybil and non-sybil identities are expected to participate are devised. A cost sharing economic mechanism called Discriminatory Rewarding Mechanism for Sybil Detection is proposed based on this game. A detective accepts a security deposit from each active agent, negotiates with the agents and offers rewards to the sybils if the latter disclose their identity. The basic objective of the detective is to determine the optimum reward amount for each sybil which will encourage the maximum possible number of sybils to reveal themselves. Maintaining privacy is an important issue for the mechanism since the participants involved in the negotiation are generally reluctant to share their private information. The mechanism has been applied to Tor by introducing a reputation scoring function.

Keywords—Game theory, Incentive mechanism, Reputation, Sybil Attack

I. INTRODUCTION

THIS explosive growth of distributed computing and universal electronic connectivity has significantly increased the chance of malicious attacks in a hostile environment. An *entity* is a software agent which may have one or more *identities* for the purpose of resource sharing, reliability and integrity (The term ‘software agent’ should not be confused with the term ‘agent’ in a game.). There are various types of attacks such as *sybil attack*, node replication attack and wormhole attack [17]. In a sybil attack a malicious entity controls multiple pseudonymous identities (sybils) and can manipulate, disrupt or corrupt an application that relies on redundancy. Messages received by an identity are communicated to its controlling entity out-of-bound [6], [10].

Asim Kumar Pal, the corresponding author, is a Professor, Management Information Systems Group, Indian Institute of Management Calcutta, Joka, Diamond Harbour Road, Kolkata, India – 700104, www.iimcal.ac.in, Tel: +919830637252 / +919748533525 (cell), +913324678300 / +913324662945, Fax: +913324678062, +913324678307, Email: asim@iimcal.ac.in, Webpage: <http://www.iimcal.ac.in/faculty/facpage.asp?ID=asim>.

Debabrata Nath is an Assistant Professor, Department of Computer Science and Engineering, Institute of Technology and Marine Engineering, Jhinga, Amira, Diamond Harbour Road, West Bengal, India – 743368, Tel: +919433170461, Email: deb3_salt@rediffmail.com.

Sumit Chakraborty was with Cognizant Technology Solutions India, in the rank of Assistant Manager, Business Development. He obtained his doctorate (Fellowship) in Management Information Systems in 2007. He has worked on problems in supply chain, negotiation, pricing and auction, and privacy preserving computation. His contacts are: Tel: +919940433441, Email: surya20046@yahoo.co.in.

To be noted that an entity is persistent, but an identity who is under the control of an entity is transient. In other words relatively speaking entity is real and identity is virtual. So when an identity gets caught or is forced to reveal itself as a sybil it can be eliminated from the system, but the controlling entity remains completely undetected and hence goes scot-free, except that some of her energy or resource is wasted because her ability to do malicious activities through the affected identity is lost. Sybil attacks may affect fair resource allocation, routing mechanisms, voting, aggregation and storage of distributed data by injecting false data or suppressing critical data [17]. Various types of distributed systems and applications are vulnerable to these attacks such as sensor and mobile ad hoc networks, auctions, reputation and trust systems, recommender systems and p2p applications [4], [11].

So far the major focus of the researchers has been on the sybil avoidance problem rather than on sybil detection. The approaches to the sybil detection are based on trusted certification [6], trusted devices, resource testing, recurring costs, direct observation [5], [7], auction scheme [8], [10], [14] and social networks [16]. Douceur’s work [6] which is a pioneering one for treating sybil attack claimed that a large-scale distributed system is highly vulnerable to sybil attacks. Margolin, Levine, et al [3], [9] – [11] have looked at the sybil attack problem and devised various incentive based methods to solve the problem. In the absence of an identification authority the conditions necessary to prevent sybil attack such as adversary’s limited capability and resource, nearly identical resource constraints for all the entities or resource constraint verification of the identities, and that all presented identities can be simultaneously validated by all entities, are not practically realizable.

Margolin and Levine [10] have proposed a sybil game based on trust and economic rationality. This scheme detects sybils by providing incentives to them for revealing their identities using an *informant* and a *detective*. The informant who pays a security deposit amount to the detective to participate in the game informs about a *target* who is a sybil under the same entity as the informant. The detective gives a reward amount to the target detected. Here, only sybils participate in the game and non-sybils have no role to play. In each round only one sybil reveals itself. There is no verification of the claim by an identity who reveals itself as a sybil. Also, since the reward amount starts from a low value (determined in a reverse auction process) in each round of the game only low-cost class sybils are detected. High-cost class sybils will reveal themselves only after low-cost class sybils are all removed from the system. The budget of the detective mainly comes from the entry fee (per identity recurring cost [9]) charged to

each identity. Cryptographic techniques were employed to protect privacy of information sharing.

The present work is based upon the scheme [10] and has proposed an improved economic game called *Sybil Detection Game* (SDG) by expanding the scope of detection of a sybil to that of by another sybil or a non-sybil, and also allowing the possibility of multiple sybil attackers. Further, a verification of the claim for a sybil is performed using a *reputation based feedback mechanism*. In a single round of the game more than one sybils can reveal themselves. High-cost class sybils, i.e. sybils with the most damaging influence, are more likely to participate in this game and reveal themselves as the rewarding mechanism is based on *negotiation*. For the case of the negotiation failing three *reward allocation mechanisms* are proposed. Several new cost parameters based on different impacts of sybil attacks on the players have been added to involve the players into rounds of negotiation which will enable them to revise their demands. This game leads to a cost sharing mechanism, namely *Discriminatory Rewarding Mechanism for Sybil Detection* (DRMSD) which attempts to maximize the number of sybils detected. In other words DRMSD minimizes the reward amount given the capacity to pay.

The above mechanism has been applied to Tor – *generation 2 Onion Routing* [18]. Tor is a distributed system where two users can communicate between them anonymously through a temporary circuit which is created by the source user. To apply the mechanism to make Tor sybil-free, the estimates of the cost parameters for the feedback are developed, besides introducing a reputation scoring function.

The structure of the paper is: Section II describes SDG with a game tree, payoffs and optimal strategies. Section III describes DRMSD along with the detailed analysis. Section IV presents the application of DRMSD to Tor. Section V concludes the paper with further scopes of research.

II. SYBIL DETECTION GAME (SDG)

Game Theory is the formal study of conflict and cooperation and is a mathematical system for analyzing and predicting how humans behave in strategic situations. Game theoretic concepts apply whenever actions of several agents are interdependent. These agents may be individuals, groups, firms or any combination of these.

In a distributed computing environment, one or more malicious entities can disrupt or corrupt an application by introducing sybil identities. One can detect those sybils by various ways. For describing such a scenario, a game theoretic approach (the SDG represented by fig. 1) is proposed. Here the sybils reveal themselves against a reward amount. The detective announces a security deposit for any identity to play the game. An identity who participates in the game is referred as a *player*, a *participant* or an *active agent*. A *distributed system administrator* monitors the overall system.

This section defines the notations required for the game (the same notations are continued throughout the paper), presents a game tree, its payoffs and optimal strategies. Also, a *tie-resolution scheme* and a feedback mechanism are described to make the game successful. Finally assumptions made for the

game are given.

A. Definitions and Notations

The following notations and definitions are used throughout the paper except that some notations are locally defined.

DSA: *distributed system administrator*.

D: *detective* recruited by DSA to detect the sybils.

$I = \{I_1, \dots, I_n\}$: n identities.

M: *mixnet* to preserve anonymity of the players from D [2].

e : *entry fee* for the identities to register, also referred as *per identity recurring cost* [9], [10].

b : *budget* amount DSA provides to D.

d : *security deposit* the identities pay to D to participate in the game.

P_1, \dots, P_m : *players* (participating identities), $m < n$.

i, j, k, l : indexes for negotiation round, player, sybil attack and correct feedback.

L_{jk} : *expected loss* of a non-sybil player from a sybil attack.

C_{jk} : *expected revenue* (or *opportunity cost*) of a sybil player from a sybil attack. This is zero for a non-sybil.

R_j : *reservation reward* demanded by a sybil against revealing itself in the current round. This is zero for a non-sybil.

R : *expected total reservation reward* by all players in the current round.

R_c : *reward capacity* of D for giving rewards.

R_f : *reward* for giving correct feedback on sybils.

P : *penalty* for wrong disclosure or feedback.

B. Game Tree

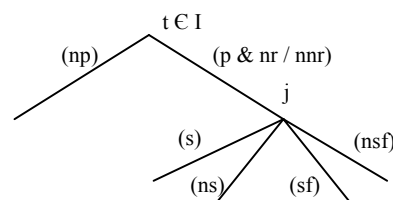


Fig. 1 Game tree of distributed system

Any identity t can choose to play (p) the game or to not play (np): a sybil identity would participate in the game to obtain a reward amount from D against revealing itself and a non-sybil identity would participate to improve the performance of the distributed system by attempting to detect sybils; otherwise they will not participate. The game will continue in one of the two modes: D will find out if negotiation is required to determine the reward amount (nr); or negotiation not required, i.e. he can directly disburse the reward amounts (nnr). Either way, a player j (either a sybil or a non-sybil) has four strategies to choose from: reveals (announces) itself as a sybil to get the reward (s); announces itself as a non-sybil to maintain reputation (ns); reveals itself as a sybil and gives feedback about other sybils to receive reward amounts (sf); and announces itself as a non-sybil and gives feedback about sybils to get reward (nsf).

C. Payoffs of the Game

The possible outcomes of the game are (np), (p, s), (p, ns), (p, sf) and (p, nsf). The payoffs of each possible outcome for the detective D and a player j (which is either a sybil S, or a non-sybil NS) are given below. The payoffs against each of the outcomes whether negotiation is required (p & nr) or not (p & nnr) remain unchanged in terms of the overall expression but not in overall values, as the values of the reward amounts change through the rounds of negotiation.

(np):	S	- [D: 0, j: $\Sigma_k C_{jk}$]
	NS	- [D: 0, j: $-\Sigma_k L_{jk}$]
(p, s):	S	- [D: $R_c - R_j$, j: $R_j + \Sigma_k C_{jk} - d$]
	NS	- [D: $R_c + P$, j: $-P - \Sigma_k L_{jk} - d$]
(p, ns):	S	- [D: $R_c + P$, j: $-P + \Sigma_k C_{jk} - d$]
	NS	- [D: R_c , j: $-\Sigma_k L_{jk} - d$]
(p, sf):	S	- [D: $R_c - R_j - \Sigma_l R_{lf}$, j: $R_j + \Sigma_k C_{jk} - d + \Sigma_l R_{lf}$]
	NS	- [D: $R_c + P - \Sigma_l R_{lf}$, j: $-P - \Sigma_k L_{jk} - d + \Sigma_l R_{lf}$]
(p, nsf):	S	- [D: $R_c + P - \Sigma_l R_{lf}$, j: $-P + \Sigma_k C_{jk} - d + \Sigma_l R_{lf}$]
	NS	- [D: $R_c - \Sigma_l R_{lf}$, j: $-\Sigma_k L_{jk} - d + \Sigma_l R_{lf}$]

D. Optimal Strategies

D always starts the game with whatever budget he receives from DSA. D does not require enter negotiation if his reward capacity exceeds the total reservation reward of the sybils. He then distributes the rewards to the sybils detected as per their demand. When there is a shortage of the reward capacity D starts the negotiation with the players. If after a given number of rounds of negotiation the situation does not improve adequately D calls optimal reward allocation mechanism [13].

An *identity* has two choices: either it will play the game or not play, as stated earlier. The optimal strategy will be to participate in the game: a sybil wants the reward and a non-sybil wants to detect the sybils.

A *sybil identity* can adopt any of the four strategies as mentioned earlier. A sybil estimates the revenue from sybil attacks it carries out. If the reward offered by D is more than expected revenue it will disclose its identity (i.e. it announces itself as a sybil); the controlling entity makes some profit as well, however at the risk of losing some of her resources, for example, the concerned identity will be removed from the system. If the reward is not adequate it may decide not to disclose its identity in which case it carries the risk of getting caught by non-sybils or other sybils and thereby adversely affects its reputation. Irrespective of these strategies a sybil may get reward for catching other sybils. A sybil also receives a penalty if it is caught giving false information about others. Thus the following can be stated about the optimal strategies of a sybil.

Theorem 1: The optimal strategy for a sybil is to disclose its identity as a sybil as well as to give feedback about other sybil identities that it comes across provided that the security deposit d is adequate satisfying equation (1) below in III.E.

A *non-sybil identity* has the same four strategies to choose

from. A non-sybil will want to prove itself as a non-sybil simply by not announcing itself as a sybil. The only incentive to declare itself as a sybil could be the reward amount. First, by doing that it takes the risk of losing reputation. If D thinks based on others' feedback that it is not a sybil then the non-sybil will be penalized. Thus the following can be stated about the optimal strategies of a non-sybil.

Theorem 2: The optimal strategy for a non-sybil is not to disclose its identity as a sybil and also to give feedback about sybil identities that it comes across provided that the security deposit d is adequate satisfying equation (1) below in III.E.

E. Tie-Resolution scheme

Suppose, a sybil x discloses its identity but before that another player y catches x . This would discourage x to participate in the game which defeats the purpose of the game. Therefore the players will be given the first chance to reveal themselves. The feedbacks will follow then. This scheme can also be used to score the reputation of the players (non-sybils) who give feedback by matching between revelations and feedbacks. Reputation building for the sybils is not required as they are eliminated from the system after receiving their reward amounts.

F. Feedback mechanism – Detection of sybils by others

Here we give an example scenario. For distributed systems such as p2p network, each identity can maintain a table of neighbours along with their reputation scores depending on the interactions between the identity and its neighbours. Each identity can treat a neighbour as a sybil based on its reputation score and can give feedback about that neighbour to D. D can then consider one as a sybil if he gets similar feedbacks from several other identities. He can also decide to launch an expensive verification mechanism such as resource testing on selected suspects. He may inform the DSA about the sybil identities detected thus or broadcast their identities. The reputation scores of the non-sybil identities who provide feedback will be updated on the basis of their predictions.

G. Assumptions

- The detective has clean reputation.
- Each agent is rational.
- If the number of players in the SDG is less than a threshold the DSA should take necessary actions to increase the number of participants, e.g. review the cost parameters.
- Privacy is an important issue in SDG so that private inputs of a player are not disclosed to others. For example, a mixnet has been proposed to keep anonymity of the players from D. This will encourage the identities to participate in the game.
- A few discrete rounds of distributed applications where each identity receives some service should be completed before D starts a fresh round of the game [8]. This helps in gathering some initial information on malicious activities performed by a few entities.

III. DRMSD

An application running in a distributed system (for examples, file sharing in a p2p system, communication among the users in a Tor network and online book store supported by a recommender system) would be infested with sybils who want to take control of a part or whole of the application. SDG is played wherein D offers rewards to the sybils if they reveal themselves. The amount of the reward can vary from one sybil to another depending on their estimates of the expected revenues from attacks, their expectation of reward amounts for revealing the truth and the reward capacity of D. The reward expectation of a sybil would be revised downward as the negotiation proceeds. The feedback about sybils by non-sybils or other sybils plays a very important role in identifying the sybils.

This section presents the algorithm of DRMSD, reward allocation mechanisms and also discusses the negotiation scenarios. Finally an analysis of the performance of DRMSD is presented and this leads to the condition required for obtaining optimal strategies for the sybils as well as the non-sybils (See Theorems 1 and 2 in Section II.D above.).

A. Algorithm

Agents: P_1, \dots, P_m, D and M .

Input: d (publicly known), b (private to D), $\sum_k C_{jk}$ & R_j (private input of each sybil j in the i -th negotiation round).

Output: The sybils learn their final reward amounts. D identifies the sybils (not necessarily all of them).

1. D announces d . m identities each pays d to D.
2. Each player j estimates and commits $\sum_k C_{jk}$ and then sends this to D through M . D computes the total (estimated) revenue of the sybils from sybil attacks and calculates his reward capacity R_c . D sets the negotiation round $i = 0$.
3. Each player j estimates and commits R_j to D through M (For a sybil $R_j \geq \sum_k C_{jk} + d$).
4. D privately computes the expected total reservation reward $R = \sum_j R_j$ and compares R with his reward capacity R_c .
5. If $R_c \geq R$, D accepts R_j for each player j as the final value of the reward. Go to step 8.
6. If $R_c < R$, D announces the start of negotiation, continues negotiation, or terminates negotiation (See sub-section C below.) as the case may be and accordingly set $i = i+1$ and go to step 3 (for revising the expected reservation rewards downward through negotiation).
7. If the negotiation finally fails, i.e. $R_c < R$, D selects optimal reward allocation mechanism (See sub-section B below.).
8. Each sybil j discloses its identity to D, who verifies these claims:
 - a. If j discloses its identity as a sybil and D finds this claim matching with the feedbacks obtained from other sybils or non-sybils he gives the value of the reward R_j to j and eliminates j from the system.
 - b. If j gives feedbacks on others which D finds to be correct from the feedbacks obtained so far, he gives reward R_f to j .

- c. If D detects any discrepancy in the disclosure or feedback of any player j , he penalizes j by P .

B. Reward Allocation Mechanisms

Three reward allocation schemes [13] are suggested here. The detective will apply this mechanism only once or more depending on the situation explained in sub-section D.

Linear reward allocation: Linear allocation is an equal sharing of the pain, i.e. shortage of reward capacity of the detective is shared among all the sybils. The reward for the sybil j is revised downward as $R_j = R_j - (R - R_c) / m'$, where m' is the number of sybils discovered so far. This is known from the number of positive reservation rewards. To be noted that some of these sybils finally may not disclose their identities (i.e. remain undetected) as their reservation reward amounts are considered too high by D.

Proportional reward allocation: The reward for the sybil j is revised downward as $R_j = R_j * R_c / R$.

Selective reward allocation: The objective is to choose the maximum possible number of sybils each receiving reservation reward R_j , i.e. the amount they had wished. The priority could be given to those who had asked for less.

C. Requirement of Negotiation

Why is the negotiation required at all? D can directly decide on the reward amounts by any reward allocation method. Negotiation is a means for the players of the game to reach mutually beneficial agreements through communication and compromise. A joint decision is made by the agents who cannot achieve their objectives through unilateral actions. They exchange information in the form of offers, counter-offers and arguments and search for a consensus. Negotiation is required with the objective to lower the reward amounts with the satisfaction of the sybils, so that the objective of detecting the maximum possible number of sybils with the minimum possible reward capacity available with D is satisfied to the maximum possible extent (See below.).

Now what will be the criterion for starting the next round of negotiation? D can announce the difference $(R - R_c)$, the total shortfall, so that the participants can reduce their reservation rewards accordingly. A sybil attacker does not know a priori when the negotiation process comes to an end and the reward allocation mechanism will be adopted, which may imply a further reduction of the reward amount given to it. To be noted is that the decision to stop the negotiation lies with D who decides based on various factors like time passed, the shortfall amount, etc. It is therefore in its (the sybil's) own interest to keep the reservation reward low. This in turn will help D to detect larger number of sybils given his reward capacity.

D. Analysis of DRMSD

Theorem 3: DRMSD utilizes the reward giving capacity of the detective to the maximum possible extent. This is in the sense that none of the sybils which remain undetected by DRMSD is *potentially able* to accept the reward offer made to

it (See below for the explanation of the term 'potentially able').

Proof: There are two strategies to detect a sybil: D detects a sybil through negotiation wherein the sybil discloses its identity, and D determines an identity as a sybil based on feedback of other sybils or non-sybils. These strategies are targeted to detect the maximum possible number of sybils. Here D may follow the tie-resolution scheme (See Section II.E above.).

D negotiates with the players (supposed to be sybils) of the game with the objective to lower the reward amounts so that more sybils can be accommodated within his reward capacity. Penalty for false claims in self disclosure of sybil identity or giving feedback about others would minimize wrong attempts or mischievous behavior.

The players would like to settle their reward claims during negotiation to avoid dispute resolution through a reward allocation mechanism (See sub-section C above.). If there is additional reward capacity of D over the demanded total reservation reward amounts this will be used as surplus for the next round of the SDG. Otherwise, the mechanism shares the shortage of reward capacity among the sybils which may possibly force a player to accept a reward amount lower than it can possibly achieve through the negotiation. D can apply his knowledge about the shortage, time passed, etc. to choose the combination or sequence of reward allocation mechanisms to optimize the negotiation process. For example, when the shortage is somewhat significant D may choose to first apply either linear or proportional reward allocation scheme to find out the players who are ready to accept the shortage. Then he can apply the selective reward allocation (an auction process) to the remaining aggrieved players (those who remain undetected even though participated in the game) and select those who can be accommodated (i.e. detected). Others will not accept the offer and remain in the system undetected (To be noted that the players communicate with D through a mixnet and hence remain anonymous.). In this sense it is a heuristic algorithm.

Now it can be claimed that DRMSD achieves the objective of utilizing the reward capacity of D to the maximum possible extent. For this the main argument is that no more aggrieved player could be accommodated, which in turn means that none of the aggrieved players could risk lowering its reservation reward. Each aggrieved player has already deposited d to participate in the game which it is surrendering now. By not accepting the reward offer it is also risking getting caught in future because of the feedback mechanism which maintains reputation score against all players. The only reason that it does not accept the offer could be that its expected gain from the sybil attack is significantly higher than the reward offer made to it. This sybil is called *potentially unable* to accept the offer made to it. Note that this notion of potential ability includes the risk perception of the individual players as well. Thus the reward capacity of D is fully utilized for detecting sybils. This completes the proof of the Theorem.

Corollary: DRMSD detects sybils having the maximum total *harming potential* each with minimum possible reward

amount as per the prevailing *market condition*. (See below for the explanation of the terms 'harming potential' of a sybil and 'market condition').

This is basically restatement of the above Theorem with the incorporation of the market condition explicit. The latter implies that the minimum possible amount of reward is disbursed to each sybil detected. This is because D would have given the maximum opportunity to the sybils for revising their reward amounts. It has been already concluded that no player remains undetected who could be accommodated for revealing itself. So for a given reward capacity D can detect the sybils having total maximum potentiality to harm. Here of course it is assumed that the harming capacity of a sybil is defined in terms of the expected gain to be achieved by it through sybil attacks.

All the above is true with the underlying assumption of market condition holding true. One has to check that how the so called market condition comes into play here. The reward capacity of D is built up through two factors: budget b given by the DSA to D and the security deposit d made by each participating player in the game. Thus the total reward capacity is $b+md$, where m is the number of players. The budget mainly comes from entry fee e obtained from each identity in the system, the surplus from previous rounds of the game, etc. The market condition involves indicators like good will of and services offered by the distributed system (e.g. Tor) and applications running on it (e.g. online selling through Tor – the purchasers remain anonymous to the seller), vulnerability of the system (e.g. sybil raters in a recommender system), etc. These market indicators indirectly determines the parameters such as b , d , e , m , m' and n (Here m' is the number of sybils discovered which can change – increase or decrease – across rounds of negotiations and n is the number of identities in the system at a particular time.). The potential for harming by a sybil is measured here through expected gains from sybil attack. The perception of this gain comes from the idea of a sybil about the vulnerability of the system, and the distribution and behaviour pattern of the users of an application. The vulnerability of the system in turn will be a complex function of the number of sybils, their behaviour patterns and also how good is the system for defending sybil attacks. Finally, the optimizing behaviour of the parties involved comes from the assumption that all of them are rational.

E. Conditions for optimal strategies

To make sure that DRMSD succeeds in its objective what is required is that b and d need to be adjusted as reward capacity R_c of D depends on these two. This is possible by adjusting e . D has to estimate d before starting the game. He does not know about m' in the system, nor he knows m . So it is difficult for D to estimate d . He must use his knowledge and experience to measure d for each round of the game.

Total of Registration fee is ne , where n is the total number of users (identities) at a particular time. Total security deposit is md , where m is the total number of players in a particular round of the game. Note, $R_c = md + b$.

Let α_j ($\alpha_j > 1$) be the *expectation factor* of a sybil over

opportunity cost to get the reservation reward, then expected total reservation reward,

$$R = \sum_j (\alpha_j \sum_k C_{jk}) + m'd.$$

Condition for fair allocation of rewards i.e., for successful negotiation is, $R_c \geq R$, i.e. $md + b \geq \sum_j (\alpha_j \sum_k C_{jk}) + m'd$, i.e.

$$d \geq (\sum_j (\alpha_j \sum_k C_{jk}) - b) / (m - m') \quad (1)$$

From the above condition followings can be said:

1. d increases with the expectation of reward
2. higher is the number of players lower should be d
3. lower the number of sybils lower can be d
4. higher is the budget lower can be d

Two scenarios of high and low entry fees are given below.

a. Low e : If e is set to a low value then the number of identities is expected to be higher which is beneficial to the system. But at the same time tendency to act as sybil will also be higher, since if the sybils are revealed (or detected by others) the controlling entity can easily (cheaply) create a new identity and go for the sybil attack again. Also for low e , the chance of failing negotiation will be higher for which D has to use the reward allocation mechanisms where dispute resolution of rewards can happen. This will discourage the sybils to play the game.

b. High e : Setting a high value of e guards the system from malicious tasks, but at the same time the number of registrations will be lower. Also, the chance of success of negotiation will be higher, or negotiation will not be required at all. This will encourage the sybils to participate in the game, because they will receive the reward amount as their expectation when revealing themselves. This will make the system sybil proof. Thus, initially even though the response is low, slowly the response will pick up. But if e is too high it may take too long to build the market.

IV. SYBIL-PROOF TOR

SDG is described earlier in the context of a p2p application (See Section II above.). Here the game in the context of Tor with the application of DRMSD is described. Tor basically is a distributed system where two users can communicate between them anonymously through a temporary circuit which is created by the source user.

This section describes Tor along with the sybil attacks on it and the application of DRMSD to Tor. The estimation of cost parameters is derived and the reputation scoring function is given for the feedback mechanism.

A. What is Tor?

The use of a switched communications network should not require revealing who is talking to whom and of course what they are talking about. *Onion Routing* is a flexible communications infrastructure that is resistant to both *eavesdropping* (to know the contents of talking) and *traffic analysis* (to know who is talking to whom). Onion Routing uses well known networking and cryptographic techniques to protect both the privacy and anonymity of Internet

communication against both eavesdropping and traffic analysis [20]. Traffic analysis can be used to infer who is talking to whom over a public network. Knowing the source and destination of Internet traffic allows others to track one's behaviour and interests. Encryption does not help against these attackers since it only hides the content of Internet traffic not the headers [18], [19].

In Tor architecture, there are several fundamental concepts which are defined as follows: An *onion router* is the server component of the network that is responsible for forwarding traffic within the core of the mix network. An *onion proxy* is the client part of the network that injects the user's traffic into the network of onion routers; one can view the onion proxy as a service that runs on the user's computer. A *circuit* is a path of three onion routers (by default) through the Tor network from the onion proxy to the desired destination. The first onion router on the circuit is referred to as the *entrance* router, the second router is called a *mix* router, and the final hop is the *exit* router. Onion proxies choose stable and high bandwidth onion routers to be entry guards which are used as an entrance router. Router information is distributed by a set of well-known and trusted *directory servers*. The unit of transmission through the Tor network is called a *cell* which is a fixed-size 512 byte packet that is padded if necessary.

Tor protects the system against traffic analysis. Onion routing proxy builds an anonymous connection through several other onion routers to the destination. The routers are also the users of the Tor system for which probably they get some incentive. Users who want to act as routers advertise their bandwidth and uptime and directory server chooses from them who have higher bandwidth and uptime. Cells are encrypted by the originator of the circuit using a layered encryption scheme. Each hop along the circuit removes a layer of encryption until the cell reaches the exit node at the end of the circuit and is fully decrypted (only the layers), reassembled into a TCP packet, and forwarded to its final destination. This layering occurs in the reverse order for data moving back to the initiator. Data passed along the anonymous connection appears different at each onion router, so data cannot be tracked en route and compromised onion routers cannot cooperate. When the connection is broken all information about the connection is cleared at each onion router [19]. A separate set of buffers are created to store cells received from each circuit. Cells are forwarded using a round-robin queuing model to give a fair amount of bandwidth to each circuit and to minimize latency.

B. Attacks on Tor

Since Tor's routing mechanism prefers high bandwidth, high uptime servers for certain portions of a flow's route, an adversary can bias the route selection algorithm toward malicious nodes with high bandwidths and high uptimes. Even adversaries with sparse resources exploit the fact that a node can lie about its resources since Tor's routing infrastructure does not verify a server's resource claims. The main aim of adversaries is to falsely advertise high bandwidth and high uptime so that they can be selected in the list of directory server and latter as entry or exit routers by others. One can

successfully register more than one router (sybil identity) all on the same IP address and different TCP port numbers so that the probability of choosing sybil routers for a user will be higher. Here the problem of adversaries with sparse resources staking claims for higher bandwidth and uptime by giving false advertisement is looked at.

C. DRMSD for Tor

Tor users (identities) first create their own circuit by selecting entry, mix and exit routers depending on high bandwidth and high uptime from the list given by directory server. After creating the circuit the source identity starts sending packets to the destination identity. An identity can observe past circuits to detect if a router's performance deviates from its advertised resource claim. This is difficult since a client cannot immediately determine which router's performance is inconsistent with its advertisement due to the multiple hop path structure.

Checking uptimes is easier than checking bandwidths. The directory server can send a simple heartbeat message periodically to test a router's uptime. This can effectively keep track how long each router has been available to the network. Directory server can then update the list of routers accordingly.

For detecting false advertised bandwidths DRMSD is used. The system has four components: authentication key generation and distribution, a sybil detection mechanism, a feedback mechanism and a reputation scoring function.

Agents: Directory Server (DS), D, M, Identities ($I_j, j=1, \dots, n$) among which DS selects routers.

Authentication key generation and distribution: Initiator of a circuit generates different communication keys for the routers of that circuit by which it can encrypt the header portion of a cell. Before layering initiator should distribute the keys to the corresponding router. Each hop along the circuit removes a layer of encryption using its own key.

Sybil detection mechanism: D performs the DRMSD in the Tor system. After completing the revelation D informs about the sybil routers to DS who then removes those routers from the list and broadcasts the information.

Feedback mechanism:

One router can be chosen by at most z identities. So the advertised bandwidth of each identity should accommodate z circuits. After selecting the routers each identity should send the router list to DS through M. If a router is selected by z identities, no more identity can further select that router. The feedback mechanism for Tor is given below.

1. Identities advertise their bandwidths and uptimes.
2. DS chooses routers between them and calculates their initial reputation scores (See sub-section E below.).
3. DS advertises List [router, bandwidth, uptime, reputation score] to the users.
4. Each sender I_j creates its circuit by selecting the routers, starts communicating with its destination and compares the actual bandwidth (p) with the minimum advertised bandwidth (q) among the routers of its circuit.
5. If $p < q$, then I_j gives feedback on the routers of its circuit along with p . These feedbacks are taken only after completing

the revelation process of DRMSD.

6. DS recalculates the reputation scores of those routers and updates the List (See sub-section E below.).

7. DS privately calculates reputation scores of the users by matching the feedbacks with the revelation. DS can use those reputation scores for considering further feedbacks.

DRMSD and feedback mechanism are run simultaneously. Most of the sybils are revealed directly from DRMSD. Sybils who have not revealed yet their reputation scores are gradually reduced according to the feedbacks. When the score of a router goes below a threshold DS can treat the former as a sybil. Also z similar feedbacks are sufficient to consider a router as a sybil. DS then removes that router from the list and informs others.

Reputation scoring function: DS generates a reputation scoring function (f_r) to calculate the reputation scores of the routers (see sub-section E below.).

D. Estimation of Cost Parameters

A possible scheme of computation of cost parameters is given here.

e : Total entry fees is ne , where n is the number of users of Tor System.

b : $b < ne$.

m : Total number of participants in the game. This is known after receiving the security deposits from the identities.

d : $d \geq (\sum_j (a_j \sum_k C_{jk}) - b) / (m - m')$. This is as per equation (1) above.

R_c : $R_c = md + b$.

L_{jk} : A function of $(p - q)$, where q = the lowest advertised bandwidth among the selected routers of a circuit, and p = actual observed bandwidth.

C_{jk} : An increasing function of $(q_j - q')$, where q_j = advertised bandwidth of a router, and q' = actual bandwidth.

R_j : $a_j \sum_k C_{jk} + d$ for a sybil, 0 for others.

R : $R = \sum_j R_j$.

R_f : Equal to d . As the revelation process is performed before the feedback mechanism, matching of these two can give a reputation of the identities. When the reputation score goes above a threshold the identities become eligible for receiving the reward amount.

P : Equal to e . Penalty should be high enough to discourage false claims.

E. Reputation Scoring Function

The reputation scoring function f_r is used by DS to maintain the reputation of the routers. Feedbacks are given only when the actual bandwidth (observed from performance) is less than the minimum bandwidth of a circuit.

Let q = minimum bandwidth of a circuit, p = actual bandwidth, initial reputation score (u) = q_j / total advertised bandwidth (this is same as the probability of the router i chosen by the users). The reputation score will be updated as, $v = u - (q - p)^2$, if $p < q$. Thus, $f_r = f(u, q, p)$.

V. CONCLUSION

In this paper a discriminatory rewarding mechanism for sybil detection has been proposed. This game follows an

economic approach with the objective that the detective detects maximum number of sybils with minimum possible reward. The computation and communication complexity of the game depends on the number of players in the game and the number of negotiation rounds. The latter in turn depends on the shortage of the reward capacity of the detective, which again depends on the expectations of the sybils. Sybils' expectation builds up from the gains these can extract from the system's vulnerability by carrying out sybil attacks. It is expected that the proposed mechanism leads to good outcomes in spite of selfish strategic behaviour by the agents. An agent can not gain incentives by misreporting its preferences to the detective. The mechanism also ensures that individual interests of the agents are best served by correct and rational behavior. The mechanism allocates rewards to the identities fairly. The reward scheme is designed in such a way that it motivates all the agents to act rationally.

In this work the number of sybil attackers is not restricted and the participating identities in the game can be both sybils and non-sybils. They can give feedback about other sybils and the detective can consider one as a sybil from a number of similar feedbacks given by different identities. Unlike [10] high-cost class sybils, i.e. sybils with most damaging influence are more likely to participate in this game in the fear that they can be detected by others. The detective verifies about the occurrence of sybil attacks through feedback so that a sybil cannot expect to receive a reward amount without initiating any attack and also by making false claims. Further a number of parameters have been used in the mechanism and their estimation is demonstrated through an application. In [10] the method uses the Dutch auction to vary the rewards. The informant protocol is used to detect one sybil in each round of execution. As the reward starts from a lower value only the low-cost class sybils are detected in every round, high-cost class sybils are not detected in any way. But the proposed mechanism can detect a number sybils belonging to both low-cost class and high-cost class sybils in a single round of the game. To be noted that the analysis of DRMSD need to be adjusted for incorporating different penalties for different kinds of violations: wrong identity disclosure by an identity and wrong feedback.

The concept of the proposed economic approach of DRMSD can be extended to various scenarios such as non-discriminatory rewarding mechanism for sybil detection and Dutch auction based mechanism for sybil detection similar to [10]. The former is suitable for a collaborative environment. This is a dynamic reward discovery mechanism where a group of agents form a coalition and negotiate with the detective. If the negotiation is successful, the value of the reward will be same for all the sybils. The latter is a hybrid mechanism having both discriminatory and nondiscriminatory reward patterns. This is a discriminatory mechanism since there are different classes of reward. The sybils corresponding to each reward class receive the same reward amount.

To make the system sybil-proof one has to correctly determine the system parameters such as entry fee, budget, and security deposit. This determination is a difficult optimization problem as outlined in the discussions on

negotiation and analysis of Section III. One may apply backward induction to do this. Further to increase the popularity of such a distributed system or application one needs to reduce the entry fee, for example. But that may encourage more sybils enter the system. What is to be done in such a situation? Further, one can look into the issue of introducing per application recurring fee. The issue of the carry over effect of feedback and reputation from one round of the game to the next needs to be explored further.

REFERENCES

- [1] A. Cheng and E. Friedman, "Sybil-proof reputation mechanisms", Proceedings of ACM SIGCOMM Workshop on Economics of Peer-to-Peer Systems, 2005, pp. 128-132.
- [2] B. Lee, C. Boyd, E. Dawson, K. Kim, J. Yang and S. Yahoo, "Providing receipt freeness in mix-net based voting protocols", Proc.- Sixth International Conference on Information Security and Cryptology, Seoul, 2003.
- [3] B.N. Levine, C. Shields and N.B. Margolin, "A Survey of Solutions to the Sybil Attack", Tech. Report, University of Massachusetts, Amherst, MA, 2006.
- [4] C. Karlof and D. Wagner, "Secure routing in wireless sensor networks: Attacks and counter measurements", Ad hoc networks Journal (Elsevier) 1 (2-3), 2003, pp. 293-315.
- [5] C. Piro, C. Shields and B. N. Levine, "Detecting the Sybil attack in adhoc networks", Proceedings of IEEE / ACM SecureComm, 2006.
- [6] J. Douceur, "The Sybil attack", Proceedings of Workshop on P2P systems (IPTPS), 2002.
- [7] J. Newsome, E. Shi, D. Song, and A. Perrig, "The Sybil attack in sensor networks: analyses & defenses", Proceedings of IPSN Intl Symposium, 2004, pp. 259-268.
- [8] M. Yokoo, Y. Sakurai and S. Matsubara, "The effect of false-name bids in combinatorial auctions: new fraud in internet auctions", Games and Economic Behavior, vol. 46, No. 1, 2004, pp. 174-188.
- [9] N.B. Margolin and B.N. Levine, "Quantifying and discouraging sybil attacks", Tech. Rep 2005-67, University of Massachusetts Amherst, 2005.
- [10] N.B. Margolin and B.N. Levine, "Informant: Detecting Sybils Using Incentives", Proceedings of Financial Cryptography, 2007.
- [11] N.B. Margolin, M. Wright and B.N. Levine, "Analysis of an incentive based protection system", Proceedings of ACM Digital Rights Management Workshop, 2004.
- [12] N. Nissan and A. Ronen, "Algorithmic mechanism design", Games Economic Behavior, 35:166-196, 2001.
- [13] S. Chakraborty, "A study of several privacy-preserving multi-party negotiation problems with applications to supply chain management" Doctoral dissertation, Indian Institute of Management Calcutta, 2007 (unpublished).
- [14] S. Rubin, M. Christodorescu, V. Ganapathy, J.T. Giffin, L. Kruger, H. Wang and N. Kidd, "An auctioning reputation system based on anomaly", Proceedings of ACM conference on Computer and Communications Security, 2005, pp. 270-279.
- [15] W. Muller, H. Plotz, J.P. Redich and T. Shiraki, "Sybil proof anonymous reputation management", SecureComm, 2008.
- [16] Y. Haifeng, M. Kaminsky, P.B. Gibbons and A. Flaxman, "Sybilguard: Defending against Sybil attacks via social networks", ACM Sigcomm'06, September 11-15, 2006, Pisa, Italy.
- [17] Y. Zhou, Y. Zhang and Y. Fang, "Access control in wireless sensor networks", Adhoc networks, Volume 5, 2007, pp. 3-13.
- [18] R. Dingleline, N. Mathewson and P. Syverson, "Tor: The second-generation onion router", In Proceedings of the 13th USENIX Security Symposium, August 2004.
- [19] I. Goldberg, "On the security of the tor authentication protocol", In Proceedings of the Sixth Workshop on Privacy Enhancing Technologies (PET 2006) (Cambridge, UK, June 2006), Springer.
- [20] J. Feigenbaum, A. Johnson, and P. Syverson, "A model of onion routing with provable anonymity", In Proceedings of the 11th Financial Cryptography and Data Security Conference (FC 2007), 2007.