# Incentivizing Responsible Networking via Introduction-Based Routing[*]

Gregory Frazier[1], Quang Duong[2], Michael P. Wellman[2], and Edward Petersen[1]

[1] BAE Systems
[2] University of Michigan

**Abstract.** The Introduction-Based Routing Protocol (IBRP) leverages implicit trust relationships and per-node discretion to create incentives to avoid associating with misbehaving network participants. Nodes exercise discretion through their policies for offering or accepting introductions. We empirically demonstrate the robustness of IBRP against different attack scenarios. We also use empirical game-theoretic techniques to assess the strategic stability of compliant policies, and find preliminary evidence that IBRP encourages the adoption of policies that limit damage from misbehaving nodes. We argue that IBRP scales to Internet-sized networks, and can be deployed as an overlay on the current Internet, requiring no modifications to applications, operating systems or core network services, thus minimizing cost of adoption.

## 1 Introduction

Militaries, governments, commercial firms, and individuals on the Internet are continually subject to a variety of attacks from a plethora of individual, organized-crime, and nation-state actors. One can identify a proximate cause for each security failure, and there exist research efforts to find mitigations to individual threat vectors as well as engineering efforts by vendors to harden products against attack. These efforts have as yet failed to significantly reduce the degree to which computer systems are being subverted.

A central cause of the widespread corruption on the Internet is the ephemeral nature of the relationship between communicating parties. Whereas there is a concrete (fiduciary, organizational) relationship between an endpoint and its Internet service provider (ISP), and between the adjacent autonomous systems that the packets traverse, the Internet protocol (IP) hides these relationships so that all two communicating endpoints know is each other's IP address. In the physical/kinetic economy, the supply chain is visible, and the resulting economic incentives reinforce good behavior. For example, when it was discovered that toothpaste exported from China contained diethylene glycol, public outrage pressured the U.S. government, which in turn pressured the Chinese government, which (then) enforced product safety laws with the exporters. Without

the exposed supply chain, consumers would have had no leverage with the exporters. Similarly, there is a supply chain behind every packet that traverses the Internet—but IP (in its current incarnation) hides it.

Introduction-based routing (IBR) seeks to limit misbehavior in a network by: i) exposing the network supply chain, ii) giving network elements authority over with whom they interact, and iii) incorporating feedback, giving participants a basis for making their relationship choices. The ability to choose with whom one interacts creates economic incentives that (in the long run) discourage network elements from interacting with nodes that repeatedly misbehave. An IBR node's decisions with regards to introductions are made by the node's *policy*, a software implementation of a strategy that maximizes the node's utility in a network that may contain corrupt participants.

The IBR protocol is completely decentralized, designed to scale to Internet-sized networks. It presumes no universal authentication mechanism, requires minimal information sharing among participants, and accommodates heterogeneous misbehavior definitions and policy implementations. Through proxy-based implementation, it can be deployed without any modifications to applications, operating systems, or the network backbone.

In the next section we describe the IBR protocol. §3 presents a simulation-based evaluation of the ability of networks of IBR-compliant nodes to create a trusted network capable of resisting various attacks. The analysis of §4 employs (empirical) game-theoretic techniques to evaluate the incentives for compliance with the protocol. In §5 we relate IBR to previous work in network security, and the final section summarizes contributions and suggests future research.

## 2   Introduction-Based Routing

Fundamental to the introduction-based approach is a node's discretion about whether to participate in an exchange of packets. Under the IBR protocol (IBRP), a packet can enter the network if and only if it is traveling on a *connection* (comparable to a VPN between two nodes). Having entered the network, it can reach only the node at the other end of the connection. Both parties to a connection must consent to participate, and either party can close the connection.

### 2.1   IBR Protocol

To establish a new connection, a node must be *introduced* by a third party with connections to both the node and the neighbor-to-be. Since no node will be connected to every other (i.e., there is no universal introducer), forming a new connection may require multiple consecutive introductions. To bootstrap participation in the network, a node must have at least one *a priori* connection (a connection that does not require an introduction). The graph of nodes linked by a priori connections defines the network.

There are three parties to an introduction: the *requester* (the node requesting the introduction); the *introducer* (the node asked to make the introduction);

and the *target* (the node to which the requester wishes to be introduced). If the introduction offer is accepted, a connection is established between the requester and target and the two nodes can exchange packets and/or request introductions to others. (These labels reflect roles in IBR, not inherent qualities of the nodes. A node may perform different roles at various points.)

A connection exists indefinitely until either party elects to close it. When a connection is closed, the requester and target provide feedback to the introducer regarding the state of the connection at the time of closing. If these nodes were introduced to the introducer, then the feedback is forwarded to those introducers after being processed. The feedback is binary—it is either positive (the other party was well-behaved) or negative. A connection closure is accompanied by messages between the two parties that signify the closure and notify the other party regarding the value of the feedback.

A key difference between IBRP and conventional routing protocols is that nodes have discretion regarding with which nodes they interact. There are two ways that a connection request can be refused. First, the introducer may respond to the introduction request with an *introduction denied* message. Second, the target may respond to the introduction offer with an *introduction declined* message. If the requester cannot find an introducer willing and able to make the introduction, then he will be unable to send packets to the target.

## 2.2   An Introduction Sequence

We illustrate the basic process of establishing a connection with a simple four-node network (Fig. 1). A priori connections exist between computers $Q$ and $M$, $M$ and $G$, and $G$ and $N$. Suppose a process on $Q$ wishes to communicate with a process on $N$. $Q$ might start by sending a request to $M$ for an introduction to $G$. [The policy on] $M$ agrees to this request, and makes the introduction. $G$ decides to accept the introduction, thus establishing a connection between $Q$ and $G$ ($\overline{QG}$). $Q$ can now ask $G$ for an introduction to $N$, which $G$ agrees to make. $N$ accepts the introduction, establishing $\overline{QN}$ and allowing the process on $Q$ to interact with the process on $N$.
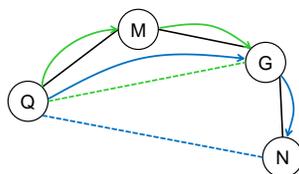


Fig. 1: Simple introduction sequence. Solid lines are a priori connections, curved arrows introductions and requests, and dashed lines the connections established.

If either $Q$ or $N$ at a later time closes the connection, both nodes provide feedback about the connection state (positive or negative) to $G$. The feedback

from $N$ is processed by $G$ and is also forwarded to $M$, since $M$ introduced the subject of the feedback ($Q$) to $G$. Nodes may play different roles on different introductions (e.g., $G$ is the target in the first introduction and the introducer in the second). Note also that both the a priori and the introduced connections are logical constructs over the network—they do not entail physical proximity between the nodes and are independent of underlying network architecture. Although the protocol does not levy any particular requirements on a node that will act as an introducer, we anticipate that in a large-scale network, introducers will typically be long-lived, well-established nodes whose capabilities as introducers are well known (possibly advertised via DNS or other lookup services).

## 2.3   Reputation

The IBR protocol does not specify the basis for nodes deciding whether to offer or accept introductions or when to close a connection. This is intentional, as we anticipate that organizations will craft policies for their nodes that correspond to their economic interests. That being said, we propose a simple model for making these decisions and argue that it i) protects the nodes' interests and globally discourages misbehavior in the network and ii) is stable, remaining an effective policy even when neighboring nodes base their decisions on a different model.

In the simple model, a node maintains a score for every neighbor to which it has a connection. We call that score a *reputation*. Much prior work in reputation-based systems (see §5) has used an explicit sharing of reputation scores as the means to curtail the attacks of a malicious participant. IBR takes a different tack. There is no sharing of reputations, and a node maintains reputation scores only for its current connections. When a target is offered an introduction, it decides whether or not to establish the connection based on the reputation of the *introducer*—not the requester, for whom the target does not have a reputation model. Therein lies the economic incentive exposed by IBR. Introducing a misbehaving node will damage the introducer's subsequent ability to provide introductions. The value of future introduction thereby drives nodes to exercise vigilance and discretion over those they are willing to introduce.

There are several reasons in practice that nodes could obtain positive value from introducing successful connections.

1. Direct benefit from activity on the connection. For instance, in a collaborative project, all participants gain from the ability of others to communicate.
2. Direct or indirect compensation from the requester for facilitating the connection. For instance, the introducer may be an ISP, who charges a fee to the requester for the introduction or more broadly for network services.
3. Implicit compensation through reciprocal action, for example the requester may be more inclined to offer introduction in kind to its existing connections.
4. Since the network itself generates value, all parties indirectly benefit from its growth to new trustworthy nodes.

We expect that factors like these provide sufficient motivation for nodes to preserve their ability to introduce, and thus cultivate and maintain a positive reputation through careful use of their introduction capability.

### 2.4  Deployment

We are pursuing the implementation of an IBR protocol proxy that is based on a bulk virtual private network (VPN) implementation. This will allow IBR to be deployed without modifications to applications, operating systems, router/switch implementations or the Internet backbone. Hosts need only to identify the IBR proxy as their gateway to participate in the protocol. Misbehavior sensors (e.g., spam detectors, network intrusion detection systems, personal security products, etc.) will be integrated with the proxy, notifying it when to preemptively close a connection and provide negative feedback.

An important element of the IBR implementation is a reliable service for discovering the introduction sequence to a given destination. We are implementing a hill-climbing algorithm, where hosts advertise nodes they accept introductions from that are higher up the hill. To obtain a connection to a given host, a node will request introductions "up the hill" until it has a connection to a node that is above the destination or has a peer that is above the destination, at which point the node requests introductions "down the hill" until it reaches the destination. This algorithm is implemented in the simulation described in the next section. In the simulation, after a warmup period where hosts establish connections to their communication partners, introductions are rarely requested. We anticipate that the same phenomena will occur in a real IBR network.

## 3  Network Performance Evaluation

Having introduced the IBR protocol, we next assess its impact on raw network performance (throughput) and its ability to discourage and/or prevent misbehavior when there are attackers present on the network.

### 3.1  Experiment Settings

We evaluate performance using an event-driven simulator. The physical details of the network were not simulated; we assumed that messages travel from source to destination in a single time unit, and we did not restrict the number of messages a node could receive per unit time. We initialized the experimental network with a set of a priori connections. Per the IBR protocol, to send a message to a given destination, a node must first establish a connection via introductions. If the connection cannot be established (due to refusal of an introducer or target server), then the message is not delivered.

Every node in an IBR network is capable of playing any of the three roles in an introduction (requester, introducer, or target). That said, for our simulation communicating nodes are explicitly designated as *clients* or *servers*. Clients generate application messages every ten to fifty time units. Each message generated is addressed to one of the server nodes in the network, selected uniformly. Upon receiving a message, servers send a response message back to the client. Nodes not designated as a client or a server provided introductions but do not participate in the application-level messaging.

Nodes in our IBR simulation monitor reputations and make introduction and connection decisions using a policy we label Compliant. A Compliant node $i$ initializes reputation scores for all a priori connections to 1.95. When $i$ accepts a connection with node $j$ through $k$'s introduction, $i$ sets its reputation score $\tau_{ij}$ for $j$ to $0.95\tau_{ik}$ where $\tau_{ik}$ is the present reputation score for the introducer $k$. When its sensors detect attacks carried out by $j$, $i$ drops $\tau_{ij}$ by 1.6 and its trust $\tau_{ik}$ for the introducer $k$ by 0.15, while the lack of sensor alerts when communicating with $j$ leads to an increase of 0.17 to $\tau_{ij}$ and of 0.1275 to $\tau_{ik}$. Upon receiving negative (positive) feedback from other nodes regarding $j$'s activities, $i$ decreases (increases) its reputation score $\tau_{ij}$ for node $j$ by 0.35 (0.01), and decrements the score of the node that introduced $j$ to $i$ by 0.05 (0.0005). During periods of no communication between $i$ and $j$, $\tau_{ij}$ regresses back toward zero with a rate of $5 \times 10^{-6} \times (t - t_{ij}^{\text{last}})^2$, where $t_{ij}^{\text{last}}$ records the last time the two nodes exchanged messages. If $j$ is an a priori connection, $i$ will disconnect from $j$ when $\tau_{ij}$ drops below 0.8. Otherwise, when $\tau_{ij}$ falls below threshold $\theta_{\text{Compliant}} = -0.4$, node $i$ disconnects from $j$ if connected, and stops introducing $j$ to others. We chose the various parameters (e.g., reputation adjustments, thresholds) through an ad hoc tuning exercise using a preliminary set of simulation runs. The strategy implemented by these constants is one of gradually increasing trust for nodes that are consistently well-behaved. When a host misbehaves, it and its introducers' reputations are significantly impacted, but it requires multiple misbehavior events in proximity before a host finds its network access restricted.

In a given experiment, one or more of the clients and/or servers are designated as *attackers*. Each message generated by an attacker node is classified as an *attack message* with some *attack probability* $q$. For the main experiments reported in §3.2, clients address attacks to servers in a uniform distribution, and attacking servers simply responded to client messages with attacks (with the specified probability). Every client and server in the network is equipped with an attack sensor that has a specified rate of detecting attack messages, as well as a false alarm (false positive) rate. Unless otherwise noted, the detection rate is 0.9 and the false positive rate 0.001 for these experiments. (These rates approximate the accuracy of spam detectors. Other forms of misbehavior and associated sensors will yield different sensitivities and specificities.)

We simulated a 4956-node network comprising 4900 client and server nodes and 56 introduction providers. The initial connection topology (Fig. 2) is a redundant tree, with seven fully connected introducers at the root and the clients and servers at the leaves. Each client and server has a single a priori connection to an introducer, notionally playing the role of ISP. There are 49 ISPs in the network, each with 99 clients and one server attached. Each ISP, in turn, is connected to two of the root introducers. Thus, each client-server connection requires between one and four introductions to establish. There is no limit to the duration of a connection or the number of connections that a node can have—once established, a connection can be used for multiple client-server transactions and is only closed when misbehavior is detected (on sensor true and false positives). We simulated conventional IP by specifying additional a priori connections from
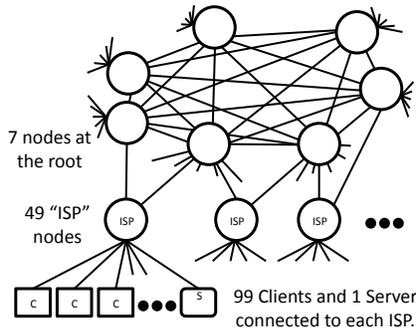
Fig. 2: The topology of a priori connections in the simulated network.

every client to every server, and employing the policy Oblivious. The threshold parameter $\theta_{\mathsf{Oblivious}} = -\infty$ means that IP nodes effectively trust all other nodes regardless of evidence about their behavior and thus never close connections. Each simulation instance generates a run that lasts for $10^4$ time units, during which every client sends and receives approximately 330 messages (3.2 million application messages per simulation). To account for the stochastic nature of the environment, for each experiment we perform 5 runs to generate one data point.

### 3.2 Results

The first experiment compared the performance of IBRP to conventional IP with a single attacking node in the network. The probability of attack $q$ was varied from 0.0 to 0.5. The results are shown in Fig. 3. Since the IP network provides no mechanism to prevent attack packets from reaching their destination, the number is directly proportional to $q$. In the IBR network, the attacker successfully sends approximately 20 attack messages. After approximately three attacks against a given server employing the Compliant policy, that server will cease to accept introductions to the attacker. After approximately five attacks, the Compliant introducer that hosts the the attacker refuses to offer any more introductions. The remaining attacks are sent on connections that the attacker already had open at the time of the fifth attack (hence the marginal drop in the number of attack messages reaching their destination for $q$ exceeding 0.2).

We note that introductions do impose some performance overhead. Under actual networking conditions, the impact is on the latency of the first packet between two hosts. Given the way we account for time in our simulation model, the introduction overhead is reflected by a loss in throughput. Each client can send only one message at a time, and the simulation is terminated after a given number of time units have elapsed. Hence the added latency of waiting to be connected reduces the total number of messages delivered. The introduction delay is experienced in the IBR network both at the beginning of the simulation (when there are no connections between clients and servers) and when a sensor
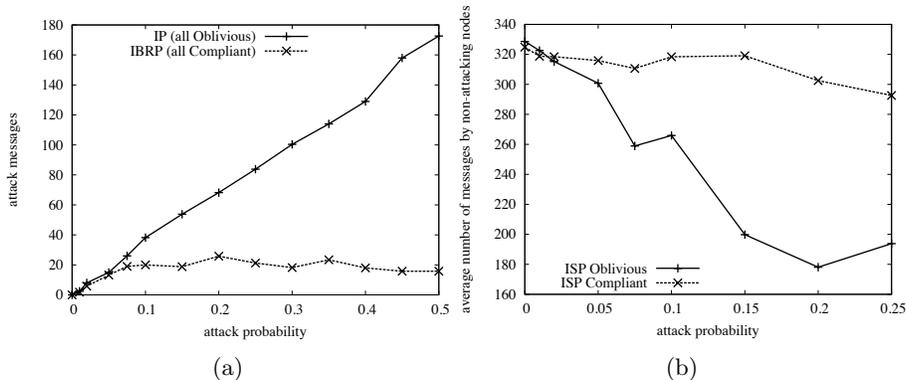
Fig. 3: (a) The number of attack messages that reached their destination per probability of attack, in IP and IBR networks of the same topology and configuration. servers. (b) Throughput experienced by the well-behaved nodes that share an ISP with the attacker, when the attacker's ISP throttles the attack (Compliant) and does not throttle the attacker (Oblivious).

false alarm causes connections to be dropped. In the IP network, a total of 3.3 million messages were delivered; in the IBR network, 3.1 million messages were delivered. (The attacker's messages were not counted.)

A tenet of IBR is that it is creating socio-economic incentives for network participants to behave responsibly. This applies not only to the network endpoints, but also to ISPs and other service providers. In the previous experiment, we showed that an ISP that applied the Compliant policy (consistent with the IBR protocol) would block the activity of an attacking node, protecting the rest of the network. However, a selfish ISP might conceivably elect not to block the attacker's activity, rather accruing subscription fees from the attacking node in return for continued service. In §4, we provide some game-theoretic evidence that rational nodes would actually adopt a compliant strategy. Here we show the performance impact on an ISP that chooses not to comply with IBR. Fig. **??** compares the performance of two networks, both of which have five attackers, all connected to the same ISP. The first network is an IBR network in which every node, including the ISP for the attackers, adopts Compliant. The second network is identical except that the attackers' ISP follows the non-compliant Oblivious policy—it continues to provide introductions for the attackers. The graph compares the throughput experienced by the *non-attacking* nodes that share the same ISP as the attackers.

When the ISP does not restrict the attackers' access to the network, all of its customers experience reduced access. This is because the second-tier introducers (the nodes to which the ISP has a priori connections) reduce the non-compliant Oblivious ISP's reputation to the point it impacts the ISP's ability to introduce

its well-behaved customers. Specifically, when negative feedback from its neighbors reduces the ISP's reputation below a threshold, the ability of the ISP's customers to reach their destinations drops precipitously. This experiment reveals a phenomenon more comprehensively analyzed in §4: that the possible benefits of providing service to a misbehaving customer will be countered by the damage caused to one's own ability to interact with the network.

### 3.3 Reputation Attacks

Of particular concern for systems that rely on reputation are attacks that undermine the reliability of information used in reputation assessment. IBR inherently mitigates the risk from such attacks by refraining from sharing reputations among nodes. Nevertheless, there remains the threat of corruption in feedback that nodes offer when closing connections. Specifically, a node might connect to its target only to close the connection and issue spurious negative feedback, hoping to convince the target's ISP to cease to introduce the target.

To assess the impact of corrupt feedback on IBR, we performed an experiment where one client node deploys such a reputation attack. Specifically, the client attacks a particular server. In our experiments, the attack rate was 0.2, which means that for every application message the client attacker generated, it had a probability of 0.2 to be addressed to the victim server. Whenever this server responds to a message from the client, the client closes the connection and sends negative feedback as though it had received an attack message from the server. For messages from all other servers, the attacker behaves like every other node, generating a false positive with probability 0.001.

IBR effectively deals with the attack in this instance. When the server discovers that the client sent negative feedback (the protocol guarantees that it will), it registers the client's misbehavior and drops its reputation. Thus, before the client can cause the server to be disconnected, the server has ceased to accept connections from the client, and the server's continued good behavior quickly restores its reputation with its ISP. In our simulations, the throughput experienced by the victim server was not adversely affected at all by the attack: the average number of messages received actually increased from 32900 to 33000.

Besides this rather mechanistic defense, the server is also shielded by a strongly positive reputation with its ISP; the attacker will have difficulty overcoming this. However, these defenses have limits. A group of coordinated attackers could achieve measurable damage to a server's reputation. This is not really surprising, as *every* protocol is susceptible to defeat by a sufficient number of Byzantine attackers [19].

### 3.4 Discussion

Our performance experiments have illustrated three key properties that IBR is designed to provide for responsible networking.

1. An ability to prevent hosts from repeatedly attacking.

2. Economic incentives for responsible network management, in the form of a reduced ability for introducers to provide introductions if they continue to introduce hosts that are known to be misbehaving.
3. Resistance to attack on based on reputation feedback.

Although the limited scope of experiments conducted reported to date prevents us from reaching sweeping conclusions, we regard the results thus far as confirming the promise of this approach.

## 4  Empirical Game Evaluation

We further evaluate whether IBRP indeed leads to a network of trustworthy connections, where malicious activity is effectively deterred and well-behaving nodes can efficiently communicate. This entails assessing to what extent IBR mechanisms induce compliant policies, in the sense that adoption of behaviors following the intended design represents a strategically stable configuration. In particular, is the use of reputation feedback enabled by IBRP for connection and introduction decisions sufficient to keep the network in a predominantly trusted state? In other words, are there plausible policies that nodes can adopt to operate both safely and efficiently? As a start toward answering these questions, we investigate the robustness and game-theoretic stability of specific compliant policies in the face of specified attacks.

### 4.1  Game Formulation

Consider an IBRP network with $n$ nodes. We construct a game model representing the interaction of these nodes over a set time horizon, as follows.

– Each node $i \in \{1, \ldots, n\}$ is a *player* (agent), who selects a *policy* (strategy) $p_i$ from a space $P_i$ of candidate policies. The policy dictates how the agent makes decisions about introductions, connections, and feedback propagation, as a function of network history.
– Nodes independently choose their policy at the beginning of a scenario run. The configuration $p = (p_1, \ldots, p_n)$ of choices is called a *policy profile*.
– The outcome of a run of the network is summarized by a *payoff* each player receives, based on how well it satisfied its communication objectives through the network over the scenario's time horizon. The payoff for node $i$ is represented by a *utility function*, $u_i : \prod_k P_k \to \Re$.

To define the game for an IBRP scenario, we specify a network topology and assign objectives for communication in the form of a distribution of messaging tasks. Receiving a benign message accrues a unit payoff, whereas a malicious message subtracts 1000 from utility. An introducer's utility is proportional to the number of successful connections it establishes. Given these scenario settings, we estimate the utility function over policy profiles, through simulations configured as described in §3.1.

We consider a policy profile strategically stable if it constitutes an (approximate) *pure-strategy Nash equilibrium* (PSNE) of the IBRP-induced game. At a Nash equilibrium, no player can increase its utility by unilaterally deviating from its policy in the profile. Given sampling noise and our rough modeling of node utility, we consider approximate degrees of strategic stability. For profile $p$, let us define player $i$'s *regret* $\epsilon_i(p)$ as the maximum gain $i$ can obtain from unilaterally changing its own policy $p_i$: $\epsilon_i(p_i, p_{-i}) = \max_{p'_i \in P_i \backslash \{p_i\}} u_i(p'_i, p_{-i}) - u_i(p_i, p_{-i})$. The overall regret of a profile $\epsilon(p) = \max_i \epsilon_i(p)$.

## 4.2 Experiment Settings

Our exploratory empirical study examines simulated small-scale IBRP networks, such as the two networks illustrated in Fig. 4. In Scenario 1 a lone introducer mediates connections between all clients and servers. Scenario 2 consists of a more elaborate network of introducers between the clients and servers.
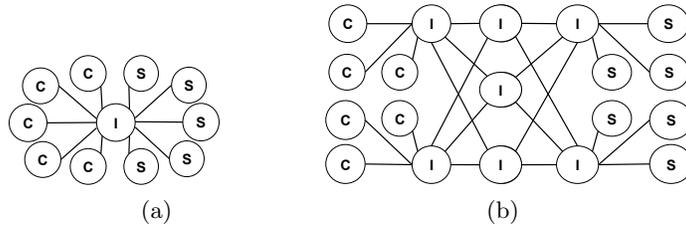


Fig. 4: (a) Scenario 1 comprises five clients and four servers, with a single introducer. (b) Scenario 2 includes a network of intermediate introducers.

We consider three candidate policies, and assume that each node chooses one and adheres to that throughout a run. The candidates include the aforementioned policy Oblivious, a parametric variant of Compliant we denote Compliant2, and a less strict version of Compliant2 named RCompliant2. We tuned Compliant2's parameters for the specific small network configurations in this section, where the presence of one attacker has a relatively more profound effect on other nodes than in larger networks. A Compliant2 node $i$ initializes all reputation values to 1.0. Each detected attack by $j$ on $i$ causes $i$ to decrease $\tau_{ij}$ by 1.0, and $\tau_{ik}$ by 0.5, where $k$ is the node that introduced $j$ to $i$. Each good transaction increases $\tau_{ij}$ by $10^{-4}$ and $\tau_{ik}$ by $5 \times 10^{-5}$. Upon receiving negative feedback regarding $j$'s activities, $i$ decreases its reputation score $\tau_{ij}$ by 0.5, and decrements the score $\tau_{ik}$ by 0.25. Node $i$ increases $\tau_{ij}$ by $5 \times 10^{-5}$ and $\tau_{ik}$ by $2.5 \times 10^{-5}$ on positive feedback about $j$. During periods of no communication between $i$ and $j$, $\tau_{ij}$ regresses back toward zero with a rate of $10^{-6}(t - t_{ij}^{\text{last}})^2$. It further specifies the threshold $\theta_{\text{Compliant2}} = -0.01$ for its neighbors' reputation. RCompliant2 retains all Compliant2's parameters except that $\theta_{\text{RCompliant2}} = -0.1$. These

candidate policies, though containing arbitrary elements, are chosen to represent a spectrum of prospective node behaviors. We are interested in comparing and contrasting the intended IBRP behaviors Compliant2 and RCompliant2 and the naive behavior Oblivious employed under IP.

For each scenario, we examine two environment settings and various levels of attack probability $q$. In the Moderate environment, nodes can perfectly detect malicious messages from their connected nodes, and only one of the client nodes is malicious. Note that since there are no attackers on the server side in Moderate, client policy is irrelevant. In environment Severe, there are attackers on both server and client sides. Moreover, nodes in Severe may mistake benign messages for harmful attacks or vice versa, each with probability 0.005.

### 4.3 Single Introducer

Under the Moderate setting of Scenario 1, for all tested values of $q$, profiles in which servers adopt either Compliant2 or RCompliant2 and the sole introducer plays Oblivious are approximate PSNEs. For the case where the attack probability is high, $q = 0.9$, the profile $p_{guard}$ in which the introducer chooses Compliant2 and the servers all play Oblivious is also a PSNE. All other profiles were found to be strategically unstable. Thus, when the sole introducer does not implement any measures to shield the servers from the attackers, the servers have to adopt compliant policies to protect themselves. In the one extreme case of reliable attack the compliant introducer with trusting servers ($p_{guard}$) achieves the same result. In all cases, the approximate equilibria dictate that either introducer or servers are compliant, with agent(s) of the other role relying on that.

Let $p_{sC}$ (resp. $p_{sRC}$) denote the approximate PSNE profile where all servers adopt Compliant2 (RCompliant2). Fig. 5 plots the difference in payoffs, denoted by $\Delta$, that a compliant server in each of these profiles would accrue by deviating to Oblivious. We also display the maximum deviation gain that a server in $p_{guard}$ would achieve by switching to a compliant policy from Oblivious. A negative value of $\Delta$ means that the agent loses utility by deviating. Note that the gains $\Delta$ are bounded above by regret, so cannot be significantly positive in approximate equilibrium. As the attack probability increases, servers in $p_{sC}$ and $p_{sRC}$ lose considerably more by deviating from their compliant policies. Fig. 5 also demonstrates that $p_{guard}$'s $\Delta$ reaches the level of approximate equilibrium only for the highest value of $q$.

Environment Severe's results are consistent with the Moderate setting in Scenario 1, for sufficiently large attack probability values. In approximate equilibrium, the non-introducer nodes adopt either Compliant2 or RCompliant2 while the introducer plays Oblivious. Let $p_{scC}$ ($p_{scRC}$) denote the profile where all non-introducer nodes adopt Compliant2 (RCompliant2). Fig. 5 confirms that when the introducer plays Oblivious, both the servers and clients have a stronger incentive to play compliant policies as $q$ increases. Smaller values of $q$ provide a less clear picture in which some approximate PSNE profiles have some non-introducer nodes playing Oblivious. Note that the possibility of false positives and negatives in detecting attacks in environment Severe, coupled with smaller values of $q$,
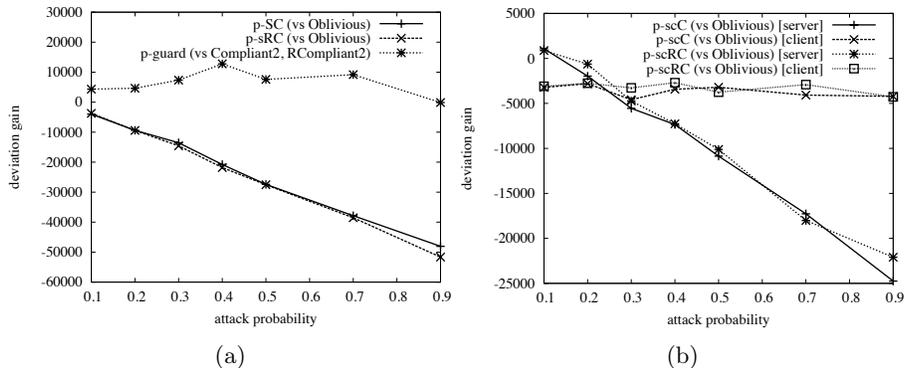
Fig. 5: (a) Servers' payoff gain in environment Moderate, when deviating to Oblivious in profiles $p_{sC}$ and $p_{sRC}$, and by deviating to either Compliant2 or RCompliant2 in $p_{guard}$. (b) Payoff gains in environment Severe, when deviating to Oblivious in profiles $p_{scC}$ and $p_{scRC}$. Results are shown for different values of the client attacker's $q$ while the server attack probability is fixed at 0.1.

makes it more difficult for nodes to effectively block potential attacks, rendering the compliant policies less desirable.

### 4.4 Network of Introducers

We conducted a further preliminary analysis of a more elaborate second scenario, with a network of introducers between clients and servers. For this scenario, we omit the RCompliant2 option, and find that all the servers adopt Compliant2 in all PSNE, under both Moderate and Severe environments. This observation reinforces the first scenario's conclusion that nodes gain substantial benefit from playing Compliant2 to protect themselves from attacks. Since a malicious client may pursue many introduction paths, it is more difficult for the network of introducers to shield the server from the need to maintain its own vigilance. While this provides appropriate incentive for the servers, it also suggests that Compliant2 policy settings may not be optimized for introducer performance. For example, note that in Scenario 2 a malicious client $C_{\mathrm{mal}}$ has multiple available paths of introduction to a given server $S_1$. When $S_1$ plays Oblivious, its proximal introducer $I_1$ essentially assumes the role of protecting $S$ from $C_{\mathrm{mal}}$'s attacks. However, given the indirect propagation of feedback, even if $I_1$ plays Compliant2 and detects the attack from one path, it may end up reintroducing $C_{\mathrm{mal}}$ to $S_1$ via a different path. A server can be introduced via three different paths that correspond to the three intermediate introducers in Scenario 2 illustrated in Fig. 4, thus allowing faster reintroductions of malicious nodes. By playing Compliant2 as well, $S_1$ uses reputation to track $I_1$'s effectiveness as a shield, and may avoid accepting the second connection from $C_{\mathrm{mal}}$.

### 4.5 Remarks

Our simulations of some simple network scenarios and environment settings, under different node policies, illustrate the possible roles of IBRP's reputation mechanism in discouraging malicious behavior. Through an empirical game-theoretic analysis, we find that the IBRP-compliant policies would be generally preferred over accepting all connection requests as in today's Internet protocol. The more hostile the environment, the greater benefit a node can obtain from adopting a compliant policy. Moreover, as others deviate by playing Oblivious, a node has greater incentive to protect itself by making reputation-based decisions. Overall, the experiments support the proposition that IBRP's reputation and feedback system enables and incentivizes nodes to adopt safer policies, thus contributing to improving network security under IBRP.

Our analysis is of course quite incomplete, and the empirical game-theoretic evidence should be regarded as supporting but preliminary. The most glaring limitation is our consideration of only three candidate policies thus far. A more comprehensive analysis would include policies that deviated from IBRP in particular ways, including selective use of reputation information, failure to propagate feedback, and false reporting of node activity. It would also include a broader palette of compliant policies, including more sophisticated ways to maintain reputation information and condition decisions on this basis. By including imperfect detection and malicious nodes, and considering a suboptimal compliant policy, we did attempt to deflect bias, but ultimately there is no substitute for the more in-depth study. This investigation should also employ much larger networks of more complex topology, and nodes with multiple roles as clients, servers, and introducers. Finally, we should also consider variations on rational behavior and their impact on IBRP performance.

## 5 Related Work

IBR participants use local trust assessments and policies to decide when to introduce and accept connections. In some respects, this resembles the "web of trust" paradigm employed notably by Pretty Good Privacy (PGP) [1, 32], in that entities rely on endorsements from trusted others in determining with whom to interact. PGP allows multiple parties to sign a given certificate, create a complex graph of trust. (PGP uses the term "introducer" to describe peer signatories.) Calculating reputation over a certificate trust-web has been extensively explored [18, 20, 26], and applied to social networks [8]. Trust-web mechanisms assume that the graph of trust relationships is generally accessible [17].

Distributed trust is used to control access in systems where centralized access control is unavailable, for example in peer-to-peer (P2P) environments [12, 13, 30]. Reputation systems [21] calculate and disseminate global reputation ratings [3, 4, 16]. Global reputation measures: i) depend on universally recognized identity and authentication for each node, ii) require uniform standards of good behavior, iii) are susceptible to attacks on reputation infrastructure, such

as Sybil and ballot-stuffing [15, 25]. TrustGuard [24] uses models of misuse to resist some of the attacks to which a shared reputation system is susceptible.

Trust management systems such as KeyNote [2] allow one to specify access control and other policies with a well-defined language. A great deal of prior work on trust and reputation mechanisms have shown how authorization decisions can be made [11, 22]. Route integrity verification mechanisms focus on improving network accountability by auditing pairwise network connections to detect malicious data-package passing activities [14, 29].

Game theory has often been used to analyze reputation systems [7], in networking and other contexts. For example, Friedman and Resnick [6] establish inherent limitations to the effectiveness of reputation systems when participants are allowed to generate new identities (a tactic IBRP is immune to), and Srivastava et al. [23] survey a range of game-theoretic applications to wireless ad hoc networks. Increasingly, game-theoretic techniques are also finding application in security domains, including but not limited to information security [9, 10, 31].

In the *empirical game-theoretic analysis* (EGTA) approach [27], expert modeling is augmented by empirical sources of knowledge, that is, data obtained through real-world observations or outcomes of high-fidelity simulation. In our own prior work, we have applied EGTA to a range of problems, including for example auctions [28] and games among privacy attackers [5].

## 6 Conclusion

The Introduction-Based Routing Protocol takes an unorthodox approach to promote a more trustworthy Internet. Rather than control access via authentication and authorization, or otherwise attempt to render forms of misbehavior algorithmically impossible, IBRP creates disincentives for associating with misbehaving nodes. Rather than attempting to track the behavior or reputation of endpoints (which is neither scalable nor practical, given the ill-defined and fleeting nature of their relationship), nodes track the (aggregate) behavior of the supply chain of network neighbors that produced the set of communication partners.

We have provided preliminary evidence that IBRP can support a trustworthy Internet. Simulations of a 5000-node network demonstrate that IBR policies successfully insulate the network from a variety of attacks. We have also shown, for some modest-sized network models, compliant policies form a strategically stable configuration, and trust decisions based on the private reputation assessment of these compliant policies limit the impact of misbehavior. We have also shown that IBRP successfully operates with imperfect sensors, with the corollary that isolated acts of misbehavior will not cause a node to be disconnected from the network.

# Bibliography

[1] Abdul-Rahman, A., Hailes, S.: A distributed trust model. In: Workshop on New Security Paradigms. pp. 48–60. Langdale, UK (1997)

[2] Blaze, M., Ioannidis, J., Keromytis, A.D.: Trust management for IPsec. ACM Transactions on Information and System Security 5(2), 95–118 (2002)

[3] Buchegger, S., Le Boudec, J.Y.: Performance analysis of the CONFIDANT protocol. In: Third International Symposium on Mobile Ad Hoc Networking and Computing. pp. 226–236. Lausanne (2002)

[4] Cornelli, F., Damiani, E., di Vimercati, S.D.C., Paraboschi, S., Samarati, P.: Choosing reputable servents in a P2P network. In: Eleventh International World Wide Web Conference. pp. 376–386. Honolulu (2002)

[5] Duong, Q., LeFevre, K., Wellman, M.P.: Strategic modeling of information sharing among data privacy attackers. Informatica 34, 151–158 (2010)

[6] Friedman, E.J., Resnick, P.: The social cost of cheap pseudonyms. Journal of Economics and Management Strategy 10(2), 173–199 (2001)

[7] Friedman, E., Resnick, P., Sami, R.: Manipulation-resistant reputation systems. In: Nisan, N., Roughgarden, T., Tardos, E., Vazirani, V.V. (eds.) Algorithmic Game Theory, pp. 677–697. Cambridge University Press (2007)

[8] Golbeck, J.A.: Computing and applying trust in web-based social networks. Ph.D. thesis, University of Maryland (2005)

[9] Grossklags, J., Christin, N., Chuang, J.: Secure or insure?: A game-theoretic analysis of information security games. In: Seventeenth International Conference on World Wide Web. pp. 209–218. Beijing (2008)

[10] Jain, M., Pita, J., Tambe, M., Ordóñez, F., Parachuri, P., Kraus, S.: Bayesian Stackelberg games and their application for security at Los Angeles International Airport. SigEcom Exchanges 7(2), 1–3 (2008)

[11] Jøsang, A., Ismail, R., Boyd, C.: A survey of trust and reputation systems for online service provision. Decision Support Systems 43(2), 618–644 (2007)

[12] Kamvar, S.D., Schlosser, M.T., Garcia-Molina, H.: The Eigentrust algorithm for reputation management in P2P networks. In: Twelfth International Conference on World Wide Web. pp. 640–651. Budapest (2003)

[13] Lagesse, B., Kumar, M., Wright, M.: AREX: An adaptive system for secure resource access in mobile P2P systems. In: Peer-to-Peer Computing'08. pp. 43–52 (2008)

[14] Laskowski, P., Chuang, J.: Network monitors and contracting systems: Competition and innovation. ACM SIGCOMM Computer Communication Review 36(4), 194 (2006)

[15] Levien, R., Aiken, A.: Attack-resistant trust metrics for public key certification. In: Seventh USENIX Security Symposium. pp. 229–42. San Antonio, TX (1998)

[16] Levine, J.: DNS Blacklists and Whitelists. RFC 5782 (Informational) (Feb 2010), http://www.ietf.org/rfc/rfc5782.txt

[17] Maurer, U.M.: Modelling a public-key infrastructure. In: Bertino, E. (ed.) Fourth European Symposium on Research in Computer Security. Lecture Notes in Computer Science, vol. 1146, pp. 325–350. Springer-Verlag (1996)

[18] Mendes, S., Huitema, C.: A new approach to the X.509 framework: Allowing a global authentication infrastructure without a global trust model. In: IEEE Symposium on Network and Distributed System Security. pp. 172– (1995)

[19] Pease, M., Shostak, R., Lamport, L.: Reaching agreement in the presence of faults. Journal of the ACM 27, 228–234 (1980)

[20] Reiter, M.K., Stubblebine, S.G.: Authentication metric analysis and design. ACM Transactions on Information System Security 2, 138–158 (1999)

[21] Resnick, P., Kuwabara, K., Zeckhauser, R., Friedman, E.: Reputation systems. Communications of the ACM 43(12), 45–48 (2000)

[22] Ruohomaa, S., Kutvonen, L.: Trust management survey. In: Third International Conference on Trust Management. pp. 77–92. Rocquencourt, France (2005)

[23] Srivastava, V., Neel, J., Mackenzie, A.B., Menon, R., DaSilva, L.A., Hicks, J.E., Reed, J.H., Gilles, R.P.: Using game theory to analyze wireless ad hoc networks. IEEE Communications Surveys and Tutorials 7(4), 46–56 (2005)

[24] Srivatsa, M., Xiong, L., Liu, L.: Trustguard: Countering vulnerabilities in reputation management for decentralized overlay networks. In: Fourteenth International Conference on World Wide Web. pp. 422–431 (2005)

[25] Sun, Y., Han, Z., Liu, K.: Defense of trust management vulnerabilities in distributed networks. IEEE Communications Magazine 46(2), 112–119 (2008)

[26] Tarah, A., Huitema, C.: Associating metrics to certification paths. In: Second European Symposium on Research in Computer Security. pp. 175–189 (1992)

[27] Wellman, M.P.: Methods for empirical game-theoretic analysis (extended abstract). In: Twenty-First National Conference on Artificial Intelligence. pp. 1552–1555. Boston (2006)

[28] Wellman, M.P., Osepayshvili, A., MacKie-Mason, J.K., Reeves, D.M.: Bidding strategies for simultaneous ascending auctions. Berkeley Electronic Journal of Theoretical Economics (Topics) 8(1) (2008)

[29] Wong, E.L., Balasubramanian, P., Alvisi, L., Gouda, M.G., Shmatiko, V.: Truth in advertising: Lightweight verification of route integrity. In: Twenty-Sixth Annual ACM Symposium on Principles of Distributed Computing. pp. 156–165. Portland, OR (2007)

[30] Xiong, L., Liu, L.: Building trust in decentralized peer-to-peer electronic communities. In: International Conference on Electronic Commerce Research (2002)

[31] Xu, J., Lee, W.: Sustaining availability of web services under distributed denial of service attacks. IEEE Transactions on Computers 52, 195–208 (2003)

[32] Zimmermann, P.R.: The Official PGP User's Guide. MIT Press (1995)