

Incentive-Compatible Interdomain Routing

Joan Feigenbaum* Christos Papadimitriou† Rahul Sami* Scott Shenker‡

January 25, 2002

Abstract

The routing of traffic between Internet domains or *Autonomous Systems* (ASs), a task known as *interdomain routing*, is currently handled by the Border Gateway Protocol (BGP). In this paper, we address the problem of interdomain routing from a mechanism-design point of view. We assume that each AS incurs a per-packet cost for carrying *transit* traffic and, in turn, is paid for carrying such traffic. The contributions of this paper are twofold. We first provide a strategyproof pricing scheme that is the only one that elicits true AS costs and does not give payments to ASs that carry no transit traffic at all. Although the transit costs are independent of the transit traffic’s source and destination, we find, somewhat counterintuitively, that the prices depend on the source and destination. Experiments with real Internet data suggest that the payments given to ASs by this scheme are unlikely to be much higher than their actual costs. We then show how this pricing mechanism can be implemented with a straightforward extension to BGP that results in only modest increases in routing-table size and convergence time. This approach of using an existing protocol as a substrate for distributed computation may prove useful in future development of Internet algorithms generally, not only for routing or pricing problems.

1 Introduction

The Internet is comprised of many separate administrative domains or *Autonomous Systems* (ASs). Routing occurs on two levels, intradomain and interdomain, implemented by two different sets of protocols. Intradomain-routing protocols, such as OSPF, route packets within a single AS. Interdomain routing, currently handled by the Border Gateway Protocol (BGP), routes packets between ASs. Although routing is a very well-studied problem, it has been approached by computer scientists primarily from an engineering or “protocol-design” perspective.

In their seminal paper on *algorithmic mechanism design*, Nisan and Ronen [NR01] advocate combining an economic, “incentive-compatibility” approach with the more traditional protocol-design approach to the problem. Internet routing is an extremely natural problem in which to consider incentives, because ownership, operation, and use by numerous independent, self-interested parties give the Internet the characteristics of an economy as well as those of a computer. In this paper, we continue the study of routing from a mechanism-design perspective, concentrating specifically on interdomain routing, for reasons explained below.

*Yale University, Computer Science Dept., New Haven, CT 06520-8285 Email: {feigenbaum, sami}@cs.yale.edu

†University of California at Berkeley, Computer Science Div., Berkeley, CA 94720. Email: christos@cs.berkeley.edu

‡ICSI, 1947 Center Street, Berkeley, CA 94704-1198. Email: shenker@icsi.berkeley.edu

We assume that each AS incurs a per-packet *cost* for carrying traffic, where the cost represents the additional load imposed on the internal AS network by this traffic. We also assume that, to compensate for these incurred costs, each AS is paid a *price* for carrying *transit* traffic, which is traffic neither originating from nor destined for that AS. Our goal is to maximize network efficiency by routing packets along the lowest-cost paths (LCPs). Standard routing protocols (such as BGP) are designed to compute LCPs given a set of AS costs. However, under many pricing schemes an AS would be better off lying about its costs;¹ such lying would cause traffic to take nonoptimal routes and thereby interfere with overall network efficiency.

To prevent this, we first ask how one can set the prices so that ASs have no incentive to lie about their costs. We require that ASs that carry no transit traffic receive no payment. We prove that there is only one *strategyproof* pricing scheme with this property; it is a member of the Vickrey-Clarke-Groves (VCG) class of mechanisms [V61, C71, G73]. We next ask how the VCG prices should be computed, and we provide a “BGP-friendly” distributed algorithm that accomplishes this.

Our results contribute in several ways to the understanding of how incentives and computation affect each other in routing-protocol design. Nisan and Ronen [NR01] and Hershberger and Suri [HS01] considered the LCP mechanism-design problem, motivated in part by the desire to include incentive issues in Internet route-selection. The LCP mechanism studied in [NR01, HS01] takes as input a biconnected graph, a single source, a single destination, and a (claimed) transmission cost for each link; the strategic agents are the links, and the mechanism computes, in a strategyproof manner, both an LCP for this single routing instance and a set of payments to the links on the LCP. This mechanism is a member of the VCG family and forms the point of departure for our work. However, our formulation of the problem differs in three respects, each of which makes the problem more representative of real-world routing:

- First, in our formulation, it is the nodes that are the strategic agents, not the links as in [NR01, HS01]. We make this choice, because we are trying to model *interdomain* routing. ASs actually *are* independent economic actors who could strategize for financial advantage in interdomain-routing decisions; in the BGP computational model into which we seek to incorporate incentive issues, it is the nodes that represent ASs and that are called upon to “advertise” their inputs to the protocol. Formulations in which the links are the strategic agents might adequately model intradomain routing, but it is not clear that incentive issues are relevant in that context; because all links and routers within a domain are owned and managed by a single entity, they are unlikely to display strategic behavior.
- Second, instead of taking as input a single source-destination pair and giving as output a single LCP, our mechanism takes in n AS numbers and constructs LCPs for all source-destination pairs. Once again, we make this choice in order to model more accurately what BGP actually does. This complicates the problem, because there are now n^2 LCP problems to solve.
- Third, we compute the routes and the payments not with a centralized algorithm, as is done in [NR01, HS01], but with a distributed protocol based on BGP. This is necessary if the motivation for the mechanism-design problem is Internet routing, because interdomain-route computation is in fact done in a distributed fashion, with the input data (AS-graph topology) and the outputs (LCPs) stored in a distributed fashion as well. The various domains are administratively separate and in some cases competitors, and there is no obvious candidate for a centralized, trusted party that could

¹Lying could increase the ASs total welfare by either attracting more traffic, and thereby increasing revenue, or increasing the price.

maintain an authoritative AS graph and tell each of the ASs which routes to use. Real-world BGP implementations could be extended easily to include our pricing mechanism, and we prove that such an extension would cause only modest increases in routing-table size and convergence time. Our approach of using an existing network protocol as a substrate for realistic distributed computations may prove useful generally in Internet-algorithm design, not only in routing or pricing problems.

Algorithm design for the Internet has the extra subtlety that adoption is not a decision by a systems manager, concerned only with performance and efficiency, but rather a careful compromise by a web of autonomous entities, each with its own interests and legacies. Backwards compatibility with an established protocol, which is one of our primary concerns here, is a constraint and criterion that is likely to become increasingly important and prevalent.

Despite these efforts to formulate the problem realistically, there are several aspects of reality that we deliberately ignore. First, it is unlikely that ASs truly incur per-packet costs, *e.g.*, in some cases transit costs are more administrative than traffic-induced. Moreover, BGP makes extensive use of policy in choosing routes, which means that LCPs are not always chosen [TGS01], and most ASs do not allow noncustomer transit traffic on their network.² In this paper, we ignore policy routing and transit restrictions; we only use LCPs. Lastly, BGP does not currently consider general path costs; it simply computes *shortest* AS paths in terms of number of AS hops. This last aspect is minor, because it would be trivial to modify BGP so that it computes LCPs; in what follows, we assume that this modification has been made.

In the next section, we provide a formal statement of the problem and in Section 3 derive the pricing scheme. In Section 4, we describe the BGP-based computational model that we use for the distributed price-calculation algorithm given in Section 5. We conclude in Section 6 with a brief discussion of open problems and future work.

2 Statement of Problem

The network has a set of nodes N , $n = \|N\|$, where each node is an AS. There is a set L of (bidirectional) links between nodes in N . We assume that this network, called the *AS graph*, is biconnected; *i.e.*, the removal of a single node does not disconnect the graph. For any two nodes $i, j \in N$, T_{ij} is the intensity of traffic (number of packets) originating from i destined for j .

We assume that a node k incurs a transit cost c_k for each transit packet it carries. For simplicity, we assume that this cost is independent of which neighbor k received the packet from and which neighbor k sends the packet to, but our approach could be extended to handle a more general case. We write c for the vector (c_1, \dots, c_n) of all transit costs and c^{-k} for the vector $(c_1, \dots, c_{k-1}, c_{k+1}, \dots, c_n)$ of all costs except c_k .

We also assume that each node k is given a payment p^k to compensate it for carrying transit traffic. In general, this payment can depend on the costs c , the traffic matrix T_{ij} , and the network topology. Our only assumption, which we invoke in Section 3, is that nodes that carry no transit traffic whatsoever receive no payment.

Our goal is to send each packet along the LCP, according to the true cost vector c . We assume the presence of a routing protocol like BGP that, given a set of node costs c , routes packets along LCPs. Furthermore, we assume that, if there are two LCPs between a particular source and destination, the

²We say that two ASs are “interconnected” if there is a traffic-carrying link between them. Interconnected ASs can be *peers*, or one can be a customer of the other. Most ASs do not accept transit traffic from peers, only from customers.

routing protocol has an appropriate way to break ties. Let $I_k(c; i, j)$ be the indicator function for the LCP from i to j ; *i.e.*, $I_k(c; i, j) = 1$, if node k is an intermediate node on the LCP from i to j , and $I_k(c; i, j) = 0$ otherwise. Note that $I_i(c; i, j) = I_j(c; i, j) = 0$; only the *transit* node costs are counted. The objective function we want to minimize is the total cost $V(c)$ of routing all packets:

$$V(c) = \sum_{i,j \in N} T_{ij} \sum_{k \in N} I_k(c; i, j) c_k$$

Minimizing V is equivalent to minimizing, for every $i, j \in N$, the cost of the path between i and j .

We treat the routing problem as a game in which the ASs are the strategic agents. Each node plays the game by reporting a transit cost. A node's transit cost is private information not known to any other node, and thus no other agent can assess the correctness of an agent's claimed transit cost. Moreover, $V(\cdot)$ is defined in terms of the true costs, whereas the routing algorithm operates on the declared costs; the only way we can be assured of minimizing $V(\cdot)$ is for agents to declare their true costs. Therefore, we must rely on the pricing scheme to incentivize agents to reveal their true costs.

To do so, we design a *mechanism* in the game-theoretic sense of the word (see [NR01]). The mechanism takes as input the AS graph and the vector c of declared costs³ and produces, as output, the set of LCPs and prices.⁴ The pricing mechanism must be *strategyproof* so that agents have no incentive to lie about their costs. For a given cost vector c , the payment p^k minus the total costs incurred by a node k is $\tau_k(c) = p^k - \sum_{i,j} T_{i,j} I_k(c; i, j) c_k$. The definition of strategyproofness is that $\tau_k(c) \geq \tau_k(c|{}^k x)$ for all x , where the expression $c|{}^k x$ means that $(c|{}^k x)_i = c_i$, for all $i \neq k$, and $(c|{}^k x)_k = x$. If this condition holds, then node k has no incentive to declare any other cost beside its true cost.

3 The Pricing Mechanism

Nisan and Ronen [NR01] described a strategyproof pricing mechanism for the LCP problem, which only involves a single source-destination pair and assumes that the links (not the nodes) are the strategic agents. Hershberger and Suri [HS01] described an efficient centralized algorithm to calculate the link prices for this single source-destination case. Our pricing mechanism must consider traffic from all source-destination pairs and must be computed in a distributed fashion; recall that both the inputs and the outputs are distributed as well, *i.e.*, that neither ever resides at a single node in the network. In this section, we derive the pricing scheme, and, in Sections 4 and 5, we describe the distributed computation.

Recall that we assume we have a biconnected graph with a routing algorithm that, when given a vector of declared costs c , will produce a set of LCPs, breaking ties in an appropriate manner; these paths are represented by the indicator functions $\{I_r(c; i, j)\}_{r \in N}$. We require that the pricing mechanism be strategyproof and that nodes that carry no transit traffic receive no payment. We now show that these two conditions uniquely determine the mechanism we must use. Moreover, we show that they require that the payments take the form of a per-packet price that depends on the source and destination; that is, the payments p^k must be expressible as

$$p^k = \sum_{i,j \in N} T_{ij} p_{ij}^k$$

³We will often use c to denote the declared costs and the true costs; usually the context will make clear which we mean.

⁴BGP will take the AS graph and c as input and produce the set of LCPs. We use this output of BGP in our mechanism and do not alter this aspect of BGP in our algorithm.

where p_{ij}^k is the per-packet price paid to node k for each transit packet it carries sent from node i destined for node j .

Theorem 1 *When routing picks lowest-cost paths and the network is biconnected, there is a unique strategyproof pricing mechanism that gives no payment to nodes that carry no transit traffic. The payments are of the form $p^k = \sum_{i,j \in N} T_{ij} p_{ij}^k$ where*

$$p_{ij}^k = c_k I_k(c; i, j) + \left[\sum_{r \in N} I_r(c|{}^k \infty; i, j) c_r - \sum_{r \in N} I_r(c; i, j) c_r \right]$$

Proof: Consider a vector of costs c . Let $u_k(c)$ denote the total costs incurred by a node for this cost vector:

$$u_k(c) = c_k \sum_{i,j \in N} T_{ij} I_k(c; i, j)$$

We can rewrite our objective function as:

$$V(c) = \sum_{i,j \in N} T_{ij} \sum_{k \in N} I_k(c; i, j) c_k = \sum_{k \in N} u_k(c)$$

Note that the routing function $\{I_k(c; i, j)\}_{k \in N}$ minimizes this quantity. The characterization of VCG mechanisms, a result due to Green and Laffont [GL79], states that the payments for any strategyproof pricing mechanism minimizing a function of the form $V(c) = \sum_{k \in N} u_k(c)$ must be expressible as

$$p^k = u_k(c) - V(c) + h_k(c^{-k})$$

where $h_k(\cdot)$ is an arbitrary function of c^{-k} . When $c_k = \infty$, we have $I_k(c|{}^k \infty; i, j) = 0$, for all i, j (because the graph is biconnected, and all other costs are finite); so (1) $p^k = 0$, by our assumption about zero payments when nodes carry no transit traffic, and (2) $u_k(c) = 0$. Thus,

$$h_k(c^{-k}) = V(c|{}^k \infty).$$

This, in turn, implies that

$$\begin{aligned} p^k &= V(c|{}^k \infty) + u_k(c) - V(c) \\ &= \sum_{i,j \in N} T_{ij} \left[c_k I_k(c; i, j) + \sum_{r \in N} I_r(c|{}^k \infty; i, j) c_r - \sum_{r \in N} I_r(c; i, j) c_r \right] \\ &= \sum_{i,j \in N} T_{ij} [p_{ij}^k], \end{aligned}$$

where

$$p_{ij}^k = c_k I_k(c; i, j) + \left[\sum_{r \in N} I_r(c|{}^k \infty; i, j) c_r - \sum_{r \in N} I_r(c; i, j) c_r \right]$$

□

This mechanism belongs to the Vickrey-Clarke-Groves (VCG) family [V61, C71, G73]. It is in essence a node-centric, all-pairs extension of the shortest-path mechanism studied by Nisan and Ronen [NR01] and Hershberger and Suri [HS01]. There are several aspects of this result that are worth noting. First,

although the payments could have taken any form and could have depended on the traffic matrix, it turns out the payments are a sum of per-packet payments that do not depend on the traffic matrix. Second, the prices $p_{i,j}^k$ are zero if the lowest cost path between i and j does not traverse k . Thus, these payments can be computed, once one knows the prices, merely by counting the packets as they enter the node. Third, although the costs did not depend on the source and destination of the packet, the prices do. Lastly, the payment to a node k for a packet from i to j is determined by the cost of the LCP and the cost of the lowest-cost path that does not path through k . We use the term *k-avoiding path* to refer to a path that does not pass through node k .

4 BGP-based Computational Model

We now seek to compute these prices $p_{i,j}^k$, using the current BGP algorithm, which is the repository of interdomain routing information, as the computational substrate. We adopt the abstract model of the BGP protocol described in [GW99], which involves several simplifying assumptions. Specifically, we assume that there is at most one link between any two ASes, that the links are bidirectional, and that each AS can be treated as an atomic entity without regard to intradomain routing issues. The network can then be modeled as a graph in which every node represents an AS, and every edge represents a bidirectional interconnection between the corresponding ASs.

BGP is a *path-vector* protocol in which every node i stores, for each AS j , the lowest-cost *AS Path* (the sequence of ASs traversed) from i to j ; in this vector, ASs are identified by their AS numbers. In addition, in our treatment, the LCP is also described by its total cost (the sum of the declared AS costs). If d is the diameter of the network (the maximum number of ASs in an LCP), a router stores $\mathcal{O}(nd)$ AS numbers and $\mathcal{O}(n)$ path costs. BGP's route computation is similar to all distance-vector routing protocols. Each router sends its routing table and, in our treatment, its declared cost, to its neighbors, and each node can then, based on this information, compute its own LCPs. When there is more than one LCP, our model of BGP selects one of them in a loop-free manner (to be defined more precisely below). As mentioned earlier, we do not consider any policy restrictions when computing these LCPs.

These routing-table exchanges only occur when a change is detected; that is, a router only sends its routing table to its neighbors when that table is different from what was sent previously. Routing tables can change either because a link was inserted or deleted (which would be detected by the nodes on either end) or when updated routing-table information is received from some other router that changes the paths (and/or costs) in the current table.⁵

The computation of a single router can be viewed as consisting of an infinite sequence of *stages*, where each stage consists of receiving routing tables from its neighbors, followed by local computation, followed (perhaps) by sending its own routing table to its neighbors (if its own routing table changed). The communication frequency is limited by the need to keep network traffic low, and hence the local computation is unlikely to be a bottleneck. Thus, we adopt as our measure of complexity the number of stages required for convergence and the total communication (in terms of the number of routing tables exchanged and the size of those tables).

If we assume that all the nodes run synchronously (exchange routing tables at the same time), BGP computes all LCPs within d stages of computation (where, again, d is the maximum number of AS hops

⁵In practice, BGP only sends the portion of the routing table that has changed. Nodes keep the routing tables received from each of their neighbors so that they can reconstruct the new routing table from the incremental update. Because the worst-case behavior is to send the entire routing table, and we care about worst-case complexity, we ignore this incremental aspect of BGP.

in an LCP). Each stage involves $\mathcal{O}(nd)$ communication on any link. The computation performed at a node i in a single stage ⁶ is $\mathcal{O}(nd \times \text{degree}(i))$.

Because this level of complexity is already deemed feasible in the current Internet, we seek to compute the prices with a similar (or better) complexity and state requirements. We describe such an algorithm in the next section.

5 Distributed Price Computation

We want to compute the p_{ij}^k using the BGP computational model described in Section 4. The input to the calculation is the cost vector c , with each c_i being known only to node i . The output is the set of set of prices, with node i knowing all the p_{ij}^k values.⁷ In describing our algorithm we assume a static environment (no route changes). The effect of removing this assumption is that the process of “converging” begins again each time a route is changed.

Our algorithm introduces additional state to the nodes and to the message exchanges between nodes, but it does not introduce any additional messages. In particular, all messages are between neighbors in the AS graph. The added state at each node consists of the reported cost of each transit node and the set of prices. This is $\mathcal{O}(nd)$ additional state, resulting in a small constant-factor increase in the state requirements of BGP. The costs and prices will be included in the routing message exchanges, and so there will be a corresponding constant-factor increase in the communication requirements of BGP.

We first investigate how the prices p_{ij}^k at node i are related to the prices at i 's neighbors.

Let $P(c; i, j)$ denote the LCP from i to j for the vector of declared costs c and let $c(i, j)$ denote the cost of this path. Define $P^{-k}(c; i, j)$ to be the lowest-cost k -avoiding path from i to j . Recall that, if there are multiple LCPs between two nodes, the routing mechanism selects one of them in a loop-free manner. *Loop-free* means that the routes are chosen so that the overall set of LCPs from every other node to j forms a tree; call this tree $T(j)$.

We treat each destination j separately. Consider the computation of p_{ij}^k for some node i at another node k on the path from i to j . Let a be a neighbor of i . There are four cases:

- **Case (i):** a is i 's parent in $T(j)$

In this case, provided that a is not k , we can extend any k -avoiding path from a to j to a k -avoiding path from i to j , and so the following inequality holds:

$$p_{ij}^k \leq p_{aj}^k \tag{1}$$

- **Case (ii):** a is i 's child in $T(j)$

Here, note that k must be on the LCP from a to j . Further, given any k -avoiding path from a to j , we can add or remove the link ia to get a k -avoiding path from i to j , and so we have:

$$p_{ij}^k \leq p_{aj}^k + c_i + c_a \tag{2}$$

⁶Because of the incremental nature of updates, where nodes need only process and forward routing entries that have changed, the communication and computational load is likely to be much less in practice.

⁷More precisely, these are the parts of the input and output that we introduce; BGP, with its standard distributed input (AS graph and costs) and distributed output (LCPs) is used as a substrate.

- **Case (iii):** a is not adjacent to i in $T(j)$, and k is on $P(c; a, j)$.

$$p_{ij}^k \leq p_{aj}^k + c_a + c(a, j) - c(i, j) \quad (3)$$

Let d be the nearest common ancestor of i and a in $T(j)$. Consider $P^{-k}(c; a, j)$, the lowest-cost k -avoiding path from a to j . We can always add the edge ia to this path to get a k -avoiding path from i to j . The inequality is then apparent by substituting the costs of the paths.

- **Case (iv):** a is not adjacent to i in $T(j)$, and k is not on $P(c; a, j)$. Let d be the nearest common ancestor of i and a . In this case, it is easy to see that

$$p_{ij}^k \leq c_k + c_a + c(a, j) - c(i, j) \quad (4)$$

Note that these four cases are not exhaustive. In particular, the cases in which a is either a descendant (but not a child) of i or $a = k$ is the parent of i are excluded. In these cases, the link ia will not be used in $P^{-k}(c; i, j)$; we can ignore neighbors in this category.

Let b be the neighbor of i on $P^{-k}(c; i, j)$; *i.e.*, the link ib is the first link on this path. We claim that, for this neighbor, the upper bounds in the previous inequalities are tight:

Lemma 1 *Let ib be the first link on $P^{-k}(c; i, j)$. Then, the corresponding inequality (1)-(4) attains equality for b .*

Proof: We can consider each of the four cases separately.

- **Case (i):**

Given that $P^{-k}(c; i, j)$ goes through its parent, it follows that b is not k , and so $p_{ij}^k = p_{bj}^k$.

- **Case(ii):** If $P^{-k}(c; i, j)$ passes through a child b , it is easy to see that $p_{ij}^k = p_{bj}^k + c_i + c_b$.
- **Case(iii):** In this case, if $P^{-k}(c; i, j)$ passes through b , it must contain $P^{-k}(c; b, j)$, and so Inequality 3 is an exact equality.
- **Case(iv):** In this case, the lowest-cost k -avoiding path through b must contain $P(c; b, j)$, and so Inequality 4 is exact.

□

Inequalities (1)-(4) and Lemma 1 together mean that p_{ij}^k is exactly equal to the minimum, over all neighbors a of i , of the right-hand side of the corresponding inequality.

Thus, we have the following distributed algorithm to compute the payment values:

The Algorithm Consider each destination j separately. The BGP table at i contains the LCP to j :

$$P(c; i, j) \equiv v_k, v_{k-1}, \dots, v_0 = j,$$

and the cost of this path, $c(i, j)$, where v_k, v_{k-1}, \dots, v_0 are the nodes on the LCP to j and $c(i, j) = \sum_{r=1}^k c_{v_r}$.

We also assume that each node knows, for each neighbor a , whether a is its parent, child, or neither in the tree $T(j)$.

At the beginning of the computation, all the entries of $p_{ij}^{v_r}$ are set to ∞ . Whenever any entry of this price array changes, the array and the path $P(c; i, j)$ are sent to all neighbors of i . As long as the network is static, the entries decrease monotonically as the computation progress. If the network is dynamic, price computation (and, as explained above, convergence) must start over whenever there is a route change.

When node i receives an updated price from a neighbor a , it performs the following updates to its internal state.

- If a is i 's parent in $T(j)$, then a scans the incoming array and updates its own values if necessary:

$$p_{ij}^{v_r} = \min(p_{ij}^{v_r}, p_{aj}^{v_r}) \quad \forall r \leq k - 1$$

- If a is a child of i in $T(j)$, i updates its payment values using

$$p_{ij}^{v_r} = \min(p_{ij}^{v_r}, p_{aj}^{v_r} + c_i + c_a) \quad \forall r \leq k$$

- If a is neither a parent nor a child, i first scans a 's updated path to find the nearest common ancestor v_t . Then i performs the following updates:

$$\begin{aligned} p_{ij}^{v_r} &= \min(p_{ij}^{v_r}, p_{aj}^{v_r} + c_a + c(a, j) - c(i, j)) & \forall r \leq t \\ p_{ij}^{v_r} &= \min(p_{ij}^{v_r}, c_k + c_a + c(a, j) - c(i, j)) & \forall r > t \end{aligned}$$

Correctness of the algorithm Inequalities (1)-(4) can be used to show by induction that the algorithm never computes a value p_{ij}^k that is too low. In order to show that the p_{ij}^k values will ultimately converge to their true values, we observe that, for every node s on $P^{-k}(c; i, j)$, the suffix of $P^{-k}(c; i, j)$ from s to j is either $P(c; s, j)$ or $P^{-k}(c; s, j)$. Thus, in the first stage of updates after the LCPs paths are found, at least one node on this path will discover its correct payment value. Inductively, we can show that all p_{ij}^k values will have their correct value after n stages.

In fact, all values will be stable after d' stages, where d' is the maximum over all i, j, k , of the number of nodes on $P^{-k}(c; i, j)$. In general, d' can be much higher than the lowest-cost diameter d of a graph. However, we don't find that to be the case for the current AS graph, as we explain in Section 6.

Using the Prices At the end of the above price computation, each node i has a full set of prices p_{ij}^k . The next question is how we can use these prices actually to compute the revenue due each node.

The simplest approach is to have each node i keep running tallies of owed charges; that is, every time a packet is sent from source i to a destination j , the counter for each node $k \neq i, j$ that lies on the LCP is incremented by p_{ij}^k . This would require $\mathcal{O}(n)$ additional storage at each node. At various intervals, nodes can send these quantities in to whatever accounting and charging mechanisms are used to enforce the pricing scheme. We assume that the submission of these running totals is done infrequently enough that the communication overhead can be easily absorbed.

In summary, we have

Theorem 2 *Our algorithm computes the VCG prices correctly, uses routing tables of size $\mathcal{O}(nd)$ (i.e., imposes only a constant-factor penalty on the BGP routing-table size), and converges in at most $(d + d')$ stages (i.e., imposes only an additive penalty of d' stages on the worst-case BGP convergence time).*

6 Conclusions and open problems

In this paper, we considered some incentive issues that arise in interdomain routing. We asked what payments are needed to elicit truthful revelation of the AS transit costs and whether they can be efficiently computed. We showed that the payments take the form of a per-packet price and that they can be computed using a simple extension to BGP that requires only a constant factor increase in communication costs. There are several promising directions for additional research.

One important issue that is not yet completely resolved is the need to reconcile the strategic model with the computational model. On the one hand, we acknowledge that ASs may have incentives to lie about costs in order to gain financial advantage, and we provide a strategyproof mechanism that removes these incentives. On the other hand, it is these very ASs that implement the distributed algorithm we have designed to compute this mechanism; even if the ASs input their true costs, what is to stop them from running a different algorithm that computes prices more favorable to them? This issue does not arise in [NR01, HS01], where the mechanism is a centralized computational device that is distinct from the strategic agents who supply the inputs, or in previous work on distributed multicast cost-sharing mechanisms [FPS01, FKSS1, FKSS2], where the mechanism is a distributed computational device (*i.e.*, a multicast tree) that is distinct from the strategic agents (who are users resident at various nodes of the tree but not in control of those nodes). If ASs are required to sign all of the messages that they send and to verify all of the messages that they receive from their neighbors, then the protocol we gave in Section 5 can be modified so that all forms of cheating are detectable [MST+01]. Achieving this goal without having to add public-key infrastructure (or any other substantial new infrastructure or computational capability) to the BGP-based computational model is the subject of ongoing further work.

There is also the issue of *overcharging*. VCG mechanisms have been criticized in the literature, because there are graphs in which the total price along a path, *i.e.*, the sum of the per-packet payments along the path, is much more than the true cost of the path. In the worst case, this total path price can be arbitrarily higher than the total path cost [AT02]. Although this is undesirable, it may be unavoidable, because VCG mechanisms are the only strategyproof pricing mechanisms for protocols that always route along LCPs. In addition, our distributed algorithm has a convergence time (measured in number of stages) d' , whereas BGP's convergence time is d ; in the worst case, $\frac{d'}{d}$ could be $\Omega(n)$. These are serious problems that could undermine the viability of the pricing scheme we present here. Thus, we ask whether these problems occur in practice. To provide a partial answer to this question, we can look at the prices that would be charged on the current AS graph if we assumed that all transit costs were the same.

Out of an 9107-node AS graph, reflecting a recent snapshot of the current Internet, we selected a 5773-node biconnected subset. We then computed d , d' , and the payments that would result from our pricing scheme, assuming a transit cost of 1 for each node. We find that $d = 8$ and $d' = 11$, and so the convergence time of the pricing algorithm is not substantially worse than that of BGP. The highest transit node price was 9, and, with uniform traffic between all pairs, the mean node payment is 1.44. In fact, 64% of the node prices were 1, and 28% of them were 2. Thus, overcharging appears not to be a problem for the current Internet.

It would be interesting to ask whether this is because of the incentive issues in AS-graph *formation*. In this paper, we merely looked at the routing aspects of a given AS graph. However, if one considers the incentives present when an AS decides whether or not to connect to another AS, the resulting transit prices would be a serious consideration. In particular, we conjecture that high node prices will not be sustainable in the Internet precisely because, if present, they would give an incentive for another AS to establish a link to capture part of that revenue, thereby driving down the transit prices. We are currently

working on models of network formation to verify this conjecture.

7 Acknowledgements

We thank Ramesh Govindan for providing us with a recent AS graph and for teaching us about the intricacies of BGP. We also thank Kunal Talwar for helpful discussions of the role of incentives in AS-graph formation.

References

- [AT02] A. Archer and E. Tardos. Frugal path mechanisms. In *Proceedings of 13th Symposium on Discrete Algorithms*, ACM Press/SIAM, New York/Philadelphia, pages 991–999, 2002.
- [C71] E. Clarke. Multipart pricing of public goods. *Public Choice* **11** (1971), pages 17–33.
- [FKSS1] J. Feigenbaum, A. Krishnamurthy, R. Sami, and S. Shenker. Approximation and Collusion in Multicast Cost Sharing, submitted. Available in preprint form at <http://www.cs.yale.edu/homes/jf/FKSS1.ps>. Abstract appears in *Proceedings of the 3rd Conference on Electronic Commerce*, ACM Press, New York, pages 253–255, 2001.
- [FKSS2] J. Feigenbaum, A. Krishnamurthy, R. Sami, and S. Shenker. Hardness Results for Multicast Cost Sharing, submitted. Available in preprint form at <http://www.cs.yale.edu/homes/jf/FKSS2.ps>.
- [FPS01] J. Feigenbaum, C. Papadimitriou, and S. Shenker. Sharing the cost of multicast transmissions. *Journal of Computer and System Sciences* **63** (2001), pages 21–41.
- [GL79] J. Green and J. Laffont. Incentives in public decision making. In **Studies in Public Economics**. Volume 1, North Holland, Amsterdam, pages 65–78, 1979.
- [GW99] T. G. Griffin and G. Wilfong. An analysis of BGP convergence properties. In *Proceedings of SIGCOMM '99*, ACM Press, New York, pages 277–288, 1999.
- [G73] T. Groves. Incentives in teams. *Econometrica* **41** (1973), pages 617–663.
- [HS01] J. Hershberger and S. Suri. Vickrey prices and shortest paths: What is an edge worth? In *Proceedings of the 42nd Symposium on the Foundations of Computer Science*, IEEE Computer Society Press, Los Alamitos, pages 129–140, 2001.
- [MST+01] J. Mitchell, R. Sami, K. Talwar, and V. Teague. Private communication, December 2001.
- [NR01] N. Nisan and A. Ronen. Algorithmic mechanism design. *Games and Economic Behavior* **35** (2001), pages 166–196.
- [TGS01] H. Tangmunarunkit, R. Govindan, and S. Shenker. Internet path inflation due to policy routing. In *Proceeding of SPIE ITCOM 2001*, SPIE Press, Bellingham, pages 19–24, 2001.
- [V61] W. Vickrey. Counterspeculation, auctions, and competitive sealed tenders. *Journal of Finance* **16** (1961), pages 8–37.