

# Congestion-Aware Internet Pricing for Media Streaming

Di Niu

Department of Electrical and Computer Engineering  
University of Alberta  
dniu@ualberta.ca

Baochun Li

Department of Electrical and Computer Engineering  
University of Toronto  
bli@eecg.toronto.edu

**Abstract**—Media webcasting and conferencing that involve many geographically distributed participants contribute significantly to congestion in the Internet. The current usage-based data pricing model does not take into account the hidden cost imposed by media streaming in the Internet core, including the network cost of replicating and relaying traffic in video multicast, and could potentially exacerbate congestion. In lieu of the recently emerged content sponsoring, in this paper, we present a simple congestion pricing model for ISPs (e.g. Comcast) to charge media streaming operators (e.g. Netflix) based on the bandwidth-delay product on each overlay link (either server-to-server or server-to-user) that the media streaming operator has chosen to use. The proposed pricing policy incentivizes different media streaming applications to collectively reduce their “waiting packets” in the Internet, alleviating congestion. We formulate the min-cost single and multiple multicast problems for the applications to construct their streaming overlays, based on a dense pool of CDN nodes. An efficient EM algorithm is given to solve the proposed geometric optimization problem and is evaluated through simulations.

## I. INTRODUCTION

Recent years have witnessed the rapid growth of Internet multimedia traffic, including video/audio streaming, Video-on-Demand (VoD), webinar/webcasting, and video conferencing, all of which are supported by applications on both personal computers and mobile devices. Most current data pricing policies place the onus on end users to pay for the consumed traffic, according to a monthly flat rate plan, a usage fee, or a combination of both. For example, a popular smartphone data plan would charge a user a monthly rate for data usage not exceeding a certain number of GBs as well as an accruing fee for each additional GB beyond the cap.

Not until recently have content providers and ISPs realized the excessive burden placed on customers in purchasing data, and due to which, their limitation or reluctance to view multimedia content or use interactive applications such as FaceTime or video streaming. In a recently initiated trend called *content sponsoring* or *application-specific pricing*, some mobile and Internet service providers have started to experiment bundling content or usage of specific applications into their data plans.

For instance, in 2011, Telus, a major ISP in Canada, offered a free six-month subscription to Rdio (music streaming) to its subscribers who purchased a Rdio-supported smartphone and data plan [1]. Similarly, a Danish operator called Tele-Danmark Communications bundled its music streaming service, TDC Play, into its mobile data plans, while Orange

France, has been bundling paid VoD, TV services and music streaming services in partnership with Deezer [2]. Similarly, in 2012, Bharti Airtel, a major Internet service provider in India, worked with Google to allow free and unlimited access to services such as Gmail, Google+ and the first pages of the websites returned by Google Search [3]. Such bundled plans are possible because many content/application providers are willing to subsidize customers by directly paying ISPs for bandwidth costs or sharing part of their advertising revenue to the ISPs. When customers use these applications, their traffic usage may not count towards the data caps in their subscribed plan or may be charged differently, encouraging customers to use the bundled multimedia applications more often.

What is the appropriate mechanism for ISPs when it comes to charging various types of content/application providers so as to better control traffic and create a positive economic cycle in the Internet? Given content/application-based pricing in its early stage, this problem is yet to be explored. In this position paper, we study how the ISPs should charge interactive media application providers including video/audio multicast operators (e.g., webcasts, video/music streaming) and multi-party conferencing operators (e.g., FaceTime, Google Hangouts) for the traffic associated with their applications. It is worth noting that unlike other forms of data flows, video/audio traffic constitutes persistent and large flows, which are potentially more harmful to the congestion condition in the Internet core.

Therefore, when pricing multimedia traffic, in addition to measuring the total amount of data transferred, it is highly desirable to impose an additional “tariff” on congested Internet routes so as to encourage application providers to select less congested routes when constructing their session-specific media streaming topology. This is analogous to charging each vehicle a higher toll on congested routes in a road network. Such a congestion-aware pricing scheme, if implemented, can help estimate and reduce the negative network externalities that each multimedia session imposes on other sessions, therefore alleviating congestion and improving the overall quality of service provided by the ISPs to all users in general.

In particular, we propose a simple congestion pricing policy, in which the ISPs charge a media streaming application provider a per-overlay-link congestion fee in proportion to its *bandwidth-delay product*. In a toy example, given two overlay routes from node  $A$  to node  $B$ , under the same

throughput, the application will choose the route with a shorter delay, thus minimizing the “waiting” packets on the overlay link, which is on expectation equal to the bandwidth-delay product according to queueing theory. This is comparable to congestion control on a road network by imposing a toll on each road according to the number of vehicles occupying the road. Under this pricing policy, the application (including single and multiple multicast sessions) will choose servers and relay nodes from the large pool of content distribution network (CDN) nodes and datacenters to form their min-cost media streaming overlay topologies, collectively minimizing the “waiting packets” occupying the Internet.

A related and important question is: does there exist an efficient and handy algorithm for any single or multiple multicast application to compute its min-cost overlay network in response to the proposed congestion pricing policy? The problem is very challenging when there is a large number of CDN nodes that can be used. To answer this question, we formulate the mathematical problem of min-cost single multicast and multiple multicast in terms of the sum of bandwidth-delay products on all the network links in a *delay space*, with the help of low-dimensional network coordinate systems [4]. We propose an efficient EM algorithm to solve this non-convex geometric optimization problem and study its performance with simulations in different scenarios.

The remainder of this paper is organized as follows. In Sec. II, we propose our congestion pricing scheme based on *bandwidth-delay products* on the links of the streaming overlay formed by an application provider. To model and study the response of multicast operators to the proposed pricing policy, we formulate the min-cost overlay multicast problem in Sec. III and propose an efficient heuristic algorithm in Sec. IV for an application to quickly form its min-cost multicast topology. In Sec. V, we further investigate how multi-party video conferencing applications should respond to the proposed congestion pricing and extend the algorithm proposed in Sec. IV to this scenario. Simulation results are presented in Sec. VI for different scenarios. Finally, we present related work in Sec. VII and conclude the paper in Sec. VIII.

## II. PRICING BASED ON BANDWIDTH-DELAY PRODUCTS

In this section, we review the existing data pricing policies targeted for end users, motivate the necessity of congestion pricing for content/application providers as a complement to end-user pricing, and describe our proposed pricing scheme inspired by some similar congestion pricing scheme in road networks.

### A. Usage-based Pricing for End Users

Current Internet or mobile service users are mostly charged in the so-called usage-based pricing model [2], [5]: users are charged either 1) a fixed monthly fee under a “Capped” plan, or 2) a “Metered” fee which is proportional to the volume of data usage, or 3) a combination of 1) and 2) in a “Cap then metered” plan, in which a user pays a flat fee up to a certain cap on data usage, beyond which the user is charged

in proportion to the usage. For example, in 2010, AT&T introduced a \$15/month data plan for 200 MB and \$25/month for 2 GB, with different rates of overage charges for the two tiers [6], [7]. Moreover, AT&T also introduced caps for its wireline service, with caps of 150 GB for DSL and 250 GB for U-Verse per month, with an overage charge of \$10 for 50 extra GB upon exceeding the caps [8]. However, it is a widespread concern that the usage-based pricing “charges customers irrespective of congestion levels in the network, and still fails to overcome the problem of large peak load costs incurred from many users crowding on the network at the same time.” [2]

### B. Congestion Pricing Inspired by Road Pricing

We do not intend to change or replace the current usage-based pricing of ISPs or mobile service providers for end users. Instead, our proposed ISP pricing policy for media streaming application providers aims to complement the pricing policy for end users by letting the ISPs charge each content/application provider a “tariff” based on the degree of congestion that the application incurs on the Internet. However, there are several challenges to designing a congestion-aware pricing policy:

- It is hard to measure the congestion level of each link: simply measuring bandwidth or delay is insufficient.
- The application may use a large number of underlying physical links in the Internet: it is impossible to monitor all the link states.
- Congestion levels change over time and are hard to be monitored dynamically, if not impossible.

Our proposed congestion-aware pricing policy is inspired by congestion-specific road pricing [2]. Transportation networks are among the first networks that adopt some form of congestion pricing, e.g., in Hong Kong [9]. One of the most natural road pricing policies is *Distance traveled pricing*, in which a vehicle pays for the distance it has traveled. *Congestion-specific pricing* [2], a dynamic pricing policy considered for Cambridge, UK, combines the distance traveled and the time spent to travel that distance; it makes the price rate per mile dependent on the speed at which the vehicle travels. As a result, each vehicle pays a fee that is proportional to *a function of both the distance traveled and its speed*.

Given a set of end users (e.g., all the users in a multicast or conferencing session), a media streaming application provider can connect the users in an overlay network by employing several CDN nodes or datacenters as content servers or relay (helper) servers. In our proposed pricing policy, we let each ISP charge the application provider a *per-minute price rate* for each link in its constructed overlay network. Such a per-minute rate is proportional to the *product of the application’s throughput on that overlay link and its packet transmission delay on the link*. As a result, the application needs to pay its ISP a per-minute congestion fee that is proportional to the *sum of bandwidth-delay products on all the overlay links* associated with the media streaming topology it has constructed.

### C. Discussions

The proposed congestion pricing policy has several advantages: *First*, the bandwidth-delay product of a certain application on a link indicates the congestion on it, while either throughput or delay alone does not. In fact, by Little’s law in queuing theory, the bandwidth-delay product is the amount of data occupying the link at any given time, i.e., the “waiting data” that has been transmitted but not yet acknowledged. If the bandwidth-delay product of the application on the link approaches the inherent bandwidth-delay product that the link can accommodate, congestion will occur. *Second*, the proposed congestion pricing is oblivious to specific underlying physical links, routers and bridges through which an application’s packets travel. By only considering the end-to-end throughput and delay on each of the overlay links among users and servers, the policy abstracts away from the detailed measurements on physical links. *Third*, it is relatively easy to record both the throughput and delay of a certain application on each overlay link by just introducing software metering functions on CDN nodes and datacenters, which may be adopted by the application as servers.

Once the congestion pricing policy is posted, every multimedia application will have the economic incentive to form a min-cost overlay network to minimize its payment to ISPs. Minimizing such a cost will then be equivalent to minimizing the sum of bandwidth-delay products on all its overlay links, or in other words, minimizing the “waiting data” it has incurred in the Internet. In the rest of the paper, we analyze how various media streaming applications should respond to the proposed congestion pricing policy. Interestingly, for both video multicast and multiparty video conferencing applications, we show that there is an algorithm for an application to effectively respond to the pricing policy, and to efficiently compute an overlay topology as well as the optimal server locations, which together minimize the operational cost of the application.

### III. MIN-COST MULTICAST IN A DELAY SPACE

Many video/audio streaming and webcasting applications can be modeled as a multicast session, where a single source node (content server) multicasts the same multimedia content to all the participating users, possibly with the help of relay servers. The application provider can freely choose the locations of relay servers from a large pool of CDN nodes and datacenters available on the Internet, and form any overlay multicast topology to connect the content server, target users and the chosen relay servers. Under the proposed congestion pricing, a min-cost multicast overlay network should be constructed to minimize the sum of bandwidth-delay products over all the overlay links in the constructed network.

It is worth noting that given the target multicast rate, the min-cost multicast problem is conventionally solved on a graph formed by the target receivers (users), the source and all the CDN nodes as candidate relay servers on the Internet to find the optimal topology and flow assignments. However, such an approach appears to be inefficient, if not impossible, when the utilizable server pool including CDN nodes and

small to medium datacenters is extremely large, leading to an overly large graph to analyze.

An alternative approach is to map the nodes onto a delay space using a *network coordinate system* [4], in which the distance between two nodes estimates their pairwise delay. As such, an “ideal” min-cost multicast network can first be computed via geometric optimization in a delay space, where we can insert relay servers at arbitrary positions in the space. Such “ideal” relay positions found in the delay space can then be mapped back to the closest real Internet servers in terms of the delay. As long as the utilizable server nodes are densely distributed, such mapping errors will be small, and this geometric optimization approach can yield a good approximation to the original min-cost multicast problem on the graph.

In this section, we formulate the min-cost multicast network in the delay space, under a constraint on the number of relay servers. Since this combinatorial problem is non-convex, we propose an effective and efficient EM algorithm to solve the problem directly in the geometric (delay) space, which can converge to local optimal solutions with high efficiency.

Although our idea can be extend to a general space, in this position paper, we focus on the min-cost multicast problem in a Euclidean space. Given  $N$  terminal nodes  $T_1, T_2, \dots, T_N$  as users with coordinates in a delay space (where the distance between two nodes models their pairwise delay on the Internet) and a multicast session from a source node  $S$  to the  $N$  terminals as sinks, the objective is to construct a min-cost network in the space, allowing the introduction of at most  $M$  extra relay nodes, and allowing any form of coding including network coding to be performed.

According to the congestion pricing proposed in Sec. II, the total congestion cost of the application is given by

$$\text{cost} = \sum_e \|e\| f(e),$$

where  $e$  is an overlay link in the constructed overlay multicast network,  $f(e)$  is the information flow rate on overlay link  $e$ , and  $\|e\|$  is the length of the link in the delay space, i.e., the end-to-end delay on overlay link  $e$ . The network cost is determined by both the positions of relay nodes and the flow assignments on the links. We call these two types of variables positions and flow assignments. Note that the flow assignments will also determine the connection topology of all the nodes, since a link exists only if there is a non-zero rate on it, and otherwise does not exist. An application provider should tune these two sets of variables with no more than  $M$  relay servers to achieve the minimum cost.

Denote  $V$  the set of all nodes and  $V_R$  the set of  $M$  candidate relay nodes, whose positions are to be found. For a node  $u \in V$ , denote  $x_u$  the position of node  $u$ . The positions of the source node and the terminal nodes are fixed input vectors, while the variables to be optimized are the positions of the relay server nodes  $V_R$  and the flow rate assignments between all pairs of nodes in the space. The geometric min-

cost multicast problem can then be formulated as

$$\underset{\substack{\{f(\vec{uv})\}_{u,v \in V}, \\ \{x_u\}_{u \in V_R}}}{\text{minimize}} \quad \sum_{u,v \in V} \|x_u - x_v\| f(\vec{uv}) \quad (1)$$

$$\text{subject to} \quad \sum_{v \in V} f_i(\vec{vu}) = \sum_{v \in V} f_i(\vec{uv}), \forall i, \forall u \in V, \quad (2)$$

$$f_i(\vec{T_i S}) = r, \quad \forall i, \quad (3)$$

$$0 \leq f_i(\vec{uv}) \leq f(\vec{uv}), \quad \forall i, \forall u, v \in V, \quad (4)$$

$$f(\vec{uv}) \leq c(\vec{uv}), \quad \forall u, v \in V, \quad (5)$$

$$|V_R| \leq M. \quad (6)$$

For every network information flow  $S \rightarrow T_i$ , there is a *conceptional flow*  $f_i(\vec{uv})$  on  $\vec{uv}$ . We call it “conceptional” because different conceptional flows can share bandwidth on the same link. Constraint (2) is the flow conservation constraint, which requires that for every node  $u \in V$ , the sum rate of all incoming conceptional flows associated with the network information flow  $S \rightarrow T_i$  equals to the sum rate of all outgoing conceptional flows associated with  $S \rightarrow T_i$ . The assigned “feedback” flow in (3) characterizes the desired receiving rate (multicast rate)  $r$  at each terminal. Constraint (5) is the trivial link capacity constraint. For every pair of nodes, we have both  $f_i(\vec{uv})$  and  $f_i(\vec{vu})$  to indicate the flows in both directions, i.e., the formed network is directed.

Constraint (4) states that the final flow rate  $f(\vec{uv})$  on any link  $\vec{uv}$  should be greater than or equal to the maximum conceptional rate among all  $i$ .  $f(\vec{uv})$  will directly affect the total cost. Finally, (6) indicates the constraint on the maximum number of relay server nodes that can be used. It has been shown [10] that in a single multicast session, any flow assignment  $\{f(\vec{uv})\}_{u,v \in V}$  satisfying (2)-(5) are feasible and can be achieved by using linear network coding (LNC) [11]. And random linear network coding (RLNC) can achieve this feasible flow assignment with high probability [12].

#### IV. AN EM ALGORITHM FOR MIN-COST MULTICAST

The real challenges to an application provider are the facts that it needs to solve problem (1) over both the variables  $x_u$  (relay positions) for  $u \in V_R$  and all the conceptional flow assignments on all the links (flow assignment) and that problem (1) is non-convex. Apparently, if an application provider cannot solve problem (1) efficiently, or in other words, derive its min-cost topology swiftly in response congestion pricing, the pricing policy proposed in Sec. II would not have achieved its expected objective — to reduce Internet congestion through fair pricing.

Fortunately, we have some good properties for problem (1), once we fix one set of variables. When the positions of relay server nodes are fixed, problem (1) is reduced to a simple linear program (LP). The number of variables is  $N + 1$  times the number of links, i.e.,  $O(N(M + N)^2)$ . The number of linear constraints is also  $O(N(M + N)^2)$ . Therefore, we can solve it efficiently with some common LP solvers.

On the other hand, when the flow assignments in the network is fixed, the cost function in (1) is the sum of

---

**Algorithm 1:** An EM algorithm for the min-cost single multicast problem (1)

---

**Input:**  $N$  terminals (users), the source node  $S$ , the constraint  $M$  on the number of relay servers.

**Output:** relay server positions, flow assignments (and topology)

---

- 1: **Initialization:** randomly generate  $M$  relay positions in the smallest box containing the source and all the terminals;

---

- 2: **Flow assignment:** fix the relay server positions, and solve the LP in (1)-(6) to get the flow rates  $f(\vec{uv})$  for all  $u, v \in V$ ;

---

- 3: **Relay position optimization:** fix the flow rates assigned in Step 2, and solve the unconstrained convex optimization problem for relay server positions  $x_u, u \in V_R$ .

---

- 4: **Random relay position regeneration:**
- 5: **for**  $i = 1$  to  $M$  **do**
- 6:   **if** the total flow on relay  $i < \epsilon$  **then**
- 7:     Regenerate a new relay  $i$  uniformly at random in the smallest box containing the source and all the terminals;
- 8:   **end if**
- 9: **end for**
- 10: Go to Step 2 unless the termination condition holds.
- 11: Remove the relay server nodes that have no flow on them.

---

- 12: **return** the relay server positions  $x_u$  for  $u \in V_R$  and flow rate assignments  $f(\vec{uv})$  for all  $u, v \in V$ .

---

norms and all the constraints (2)-6) are irrelevant to the relay positions. Problem (1) now reduces to a convex optimization problem. There are many efficient algorithms to solve such kind of problems. More specifically, for the sum of norms in this problem, an *Equilibrium method* based on the notion of forces and stretching has been proposed in [13], which can efficiently converge to the optimal solution to relay server positions.

We propose an EM heuristic algorithm for application providers to solve problem (1), based on the two observations above. In the EM algorithm, the above two local optimizations for relay positions and flow assignments are alternately performed until convergence. Our proposed EM algorithm is shown in Algorithm 1. Initially, Step 1 randomly assigns the positions of the relay server nodes in the smallest box region containing all the terminals and the source. The following steps are iterative operations. In each iteration, there are three major steps. We first solve the LP in (1)-(6) with the relay positions fixed to obtain the flow rate assignments. Then with these flow rates fixed, we solve a convex optimization problem for the relay server node positions. Furthermore, for each relay node that has no throughput on it, we randomly reassign a new position to it and repeat the iterations. Note that  $\epsilon$  is a small positive threshold to exclude the fake non-zeros, since in our LP solver, there is always a small non-zero value on an actually zero-valued variable.

We introduce a counter to help us monitor the termination condition. In each iteration, we first calculate the change of the cost from the previous iteration, and increment the counter

if the changed ratio is less than a certain threshold. Once the counter reaches some number, we terminate the iterations and remove the relay nodes that have no throughput.

In reality, each application can run Algorithm 1 to obtain the ideal relay server positions  $x_u$  for  $u \in V_R$  and flow rate assignments  $f(\vec{uv})$  for all  $u, v \in V$ . Then, each ideal relay server position can be mapped to the *closest* real physical server in the candidate server pool of CDN nodes and datacenters.

## V. MULTIPARTY VIDEO CONFERENCING

In (multiparty) video conferencing applications such as Google Hangouts, FaceTime and Skype, each participating user streams its locally captured video to all other users. We now study how these applications should construct its media streaming overlay given the locations of the end users.

Given a set of  $N$  participating terminals, a multiparty conferencing session can be modelled as a special case of multiple multicast, where each terminal serves as a source with all the other nodes being its receivers. In general, in a multiple multicast session, denote the sources as  $S_j \in V, j = 1, \dots, N$ . Each source  $S_j$  has a multicast rate  $r_j$  and multiple target receivers  $T_{ij} \in V, i = 1, \dots, k_j$ . Note that a node  $u$  can be a receiver of more than one source. Consider the session associated with each source  $S_j$  as an *independent* single multicast, with conceptual flows  $\{f_{ij}(\vec{uv}) | i = 1, \dots, k_j\}$ , where  $f_{ij}(\vec{uv})$  is the flow rate associated with the network information flow  $S_j \rightarrow T_{ij}$  on link  $\vec{uv}$ .

Now the min-cost multiple multicast overlay network should be constructed by solving the following problem:

$$\underset{\substack{\{f(\vec{uv}) | u, v \in V\}, \\ \{x_u | u \in V_R\}}}{\text{minimize}} \quad \sum_{u, v \in V} \|x_u - x_v\| f(\vec{uv}) \quad (7)$$

$$\text{subject to} \quad \sum_{v \in V} f_{ij}(\vec{vu}) = \sum_{v \in V} f_{ij}(\vec{uv}), \forall i, \forall j, \forall u \in V, \quad (8)$$

$$f_{ij}(\vec{T}_{ij} S_j) = r_j, \quad \forall i, \forall j, \quad (9)$$

$$\sum_{j=1}^N \max_{i | T_{ij} S_j \neq \vec{uv}} \{f_{ij}(\vec{uv})\} \leq f(\vec{uv}), \forall u, v \in V, \quad (10)$$

$$f(\vec{uv}) \leq c(\vec{uv}), \quad \forall u, v \in V, \quad (11)$$

$$f_{ij}(\vec{uv}) \geq 0, \quad \forall i, \forall j, \forall u, v \in V, \quad (12)$$

$$|V_R| \leq M. \quad (13)$$

Constraint (8) is the flow conservation rule for every network information flow  $S_j \rightarrow T_i$ . The assigned “feedback” flow in (9) characterizes the desired multicast rate  $r_j$  for each source  $S_j$  to its receiver set  $\{T_{ij} | i = 1, \dots, k_j\}$ . Constraint (10) gives the final flow rate  $f(\vec{uv})$  on each link  $\vec{uv}$ . To see the rationale, now consider the single multicast associated with each source  $S_j$ . Allowing *intra-session* linear network coding, the final flow rate associated with source  $S_j$  on link  $\vec{uv}$  should be at least  $\max_i \{f_{ij}(\vec{uv})\}$  (excluding the “feedback” flows) for the single multicast session to be feasible, according to [10]. If we do not allow *inter-session* network coding and treat

the  $N$  sources in  $N$  independent single multicast sessions, then the final flow rate  $f(\vec{uv})$  assigned on each link  $\vec{uv}$  should be at least the sum of final flow rates associated with all the sources  $S_j$  on link  $\vec{uv}$ , explaining constraint (10). Assuming that no *inter-session* network coding is performed, constraints (8)-(13) outline the largest region of all feasible flow assignments in the given multiple multicast session. Just like problem (1), the min-cost multiple multicast problem (7) can also be solved using an EM algorithm similar to Algorithm 1. Once the ideal relay server positions are found, they can be mapped to the closest physical servers in the pool of CDN nodes and datacenters.

## VI. SIMULATIONS

To evaluate the effectiveness of the proposed pricing scheme and the responses of media streaming application providers, we simulate three different application scenarios, including two single multicast sessions and one video conferencing session, shown in Fig. 1. For each session, we assume that the coordinates of all the terminals (participating users) can be obtained in a 2-D delay space before the session starts by probing some static server nodes via ping [4].

For a 10-user single multicast with source rate  $r = 1$ , Fig. 1(a) shows the min-cost overlay topology with relay server positions in the delay space computed by Algorithm 1. The resulted cost in terms of sum of bandwidth-delay products on all links is 2.31. In comparison, the cost of direct sending, currently a common solution where the source server sends a separate stream to each terminal at a rate of 1 without adopting relay CDN nodes, is 5.24, which more than doubles the cost of Algorithm 1. Fig. 1(b) shows the min-cost topology computed for another multicast session, where terminals are distributed differently. In this case, the topology computed by Algorithm 1 must use network coding to achieve the desired multicast rate 1. The cost is 5.15, as compared to 9.7 for direct sending.

Fig. 1(c) plots the min-cost overlay topology with relay server positions computed for a 5-user video conferencing session, which is multiple multicast session with each terminal being a source for all other terminals. The total resulted cost is 20.94, as compared to 28.4 for direct sending (P2P mode without server) and 23.05 for using a single relay server at the centroid of the terminals. In all three cases, the cost reduction in terms sum of bandwidth-delay products is significant.

## VII. RELATED WORK

The concept of responsive pricing for congestion control has existed for a long time. In close-loop feedback pricing [14], the network load, measured in terms of buffer occupancy at the gateway, is converted to a price per packet for users’ adaptive applications to decide how much data to transmit. In a study of revenue and welfare maximization for customer calls [15], users initiate calls that have different resource requirements and call duration. Based on the network congestion level, the service provider charges a fee per call, which in turn affects the user demand. Time-dependent usage-based pricing [16] assumes some form of utility functions adopted by customers,

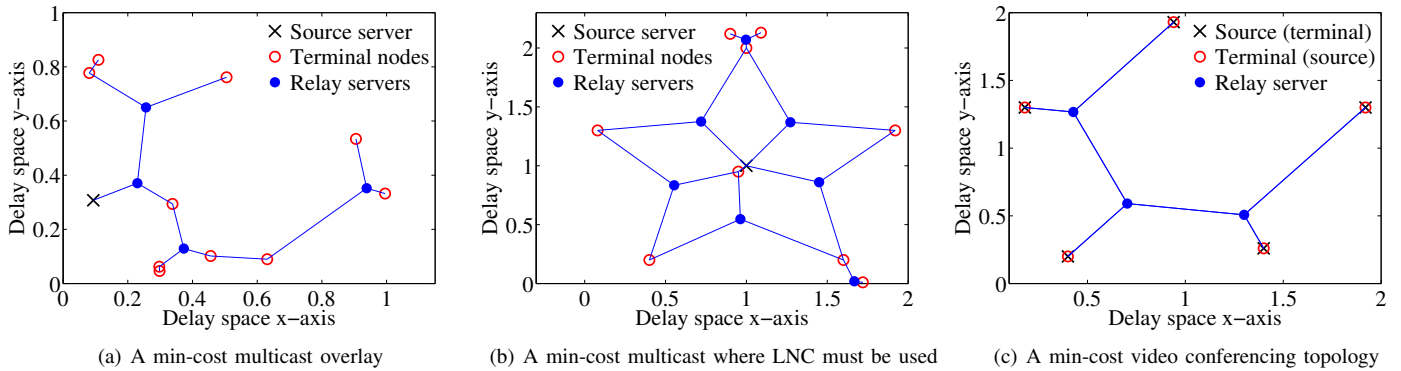


Fig. 1. The min-cost overlay topologies of three different sessions in a 2-D delay space with source rates all equal to 1. Terminal nodes are end users, servers can be chosen from the pool of CDN nodes and datacenters. In the conferencing session, each terminal is also a source to other terminals.

and aims to compute the dynamic prices to be offered to customers, using convex optimization, with the objective of minimizing the cost of overusing capacity on bottleneck links and shifting away peak demand. This paper is similar to the above work in that pricing is dependent on network states. However, instead of using pricing to influence or shift user demands [14]–[16], we use pricing to indirectly control the way that application providers construct their session-specific overlay structures on which to route their traffic. Under our proposed pricing policy based on bandwidth-delay products, every application will be incentivized to minimize its aggregate “waiting data” incurred on the Internet, thus alleviating congestion.

The geometric min-cost multicast problem (1) in the delay space is similar to the *Space Information Flow* (SIF) problem [13], which aims to minimize the sum of bandwidth-distance products in a (e.g., geographic) space, allowing network coding and free insertion of relay nodes. [13] presents a heuristic solution to the space information flow problem. However, it cannot solve the problem under a relay number constraint. Also, the solution in [13] approximates the geometric problem with a graph version of the problem by dividing the space with fine-grained grids. As a large number of intermediate variables are introduced, it has a large complexity when the dimension of the space is high. Our EM algorithm is a heuristic solution to SIF under a relay server number constraint, and yields fast convergence. The performance optimization of multiparty conferencing in terms of other metrics such as the sum of end-to-end delays has been studied in [17].

## VIII. CONCLUDING REMARKS

In this paper, we have proposed a congestion pricing policy for media streaming application providers including video/audio multicast and multiparty conferencing. Assuming the applications can freely employ CDN nodes and datacenters to form their streaming overlay, we let ISPs charge an application provider a fee in proportion to the sum of bandwidth-delay products on its overlay links, in order to encourage the applications to form low-cost single or multiple multicast overlays. Since such defined cost is directly related to the amount of “waiting packets” occupying the Internet, in response to the pricing scheme, the participating applications

will collectively reduce the congestion levels they incur on the Internet. We provide efficient heuristic solutions to min-cost multicast problems and show through simulations that an optimized overlay topology employing additional server nodes can cut the cost down for nearly one half as compared to the most common practice of direct sending. This implies that once the proposed pricing policy is adopted, media streaming applications will have a strong incentive to optimize their cost, thus indirectly alleviate the Internet congestion conditions.

## REFERENCES

- [1] M. Vardy, “Telus tuneage: Now powered by rdio,” *The Next Web*, <http://thenextweb.com/ca/2011/08/03/telus-tuneage-now-powered-by-rdio/>.
- [2] S. Sen, C. Joe-Wong, S. Ha, and M. Chiang, “Pricing data: A look at past proposals, current plans, and future trends,” Tech. Rep., 2012.
- [3] P. Goldstein, “Google joins with India’s Bharti Airtel for toll-free wireless Internet service,” <http://www.fiercemobileit.com/story/att-ceo-content-providers-asking-toll-free-data-plans/2012-07-18>.
- [4] F. Dabek, R. Cox, F. Kaashoek, and R. Morris, “Vivaldi: a decentralized network coordinate system,” in *Proc. ACM SIGCOMM*, 2004.
- [5] J. Walrand, *Economic Models of Communication Networks*. Springer Publishing Company, 2008.
- [6] D. Frakes, “AT&T announces tethering details and new plans for iPhone, iPad,” *InfoWorld*, June 2 2010.
- [7] C. Kang, “AT&T wireless scraps flat-rate internet plan,” *The Washington Post*, 2010.
- [8] P. Taylor, “AT&T imposes usage caps on fixed-line broadband,” *Financial Times*, March 14 2011.
- [9] T. D. Hau, “Electronic road pricing: Developments in hong kong 1983-1989,” *Journal of Transport Economics and Policy*, vol. 24, no. 2, pp. 203–214, May 1990.
- [10] Z. Li, “Min-cost multicast of selfish information flows,” in *Proc. IEEE INFOCOM*, 2007.
- [11] S. R. Li, R. W. Yeung, and N. Cai, “Linear network coding,” *IEEE Transactions on Information Theory*, vol. 49, p. 371, 2003.
- [12] T. Ho, M. Medard, J. Shi, M. Effros, and D. Karger, “On randomized network coding,” in *Proc. of Allerton Conference*, 2003.
- [13] J. Huang, X. Yin, X. Zhang, X. Du, and Z. Li, “On space information flow: Single multicast,” in *Proc. the International Symposium on Network Coding (NetCod)*, 2013.
- [14] J. Murphy and L. Murphy, “Bandwidth allocation by pricing in ATM networks,” *IFIP Trans. C: Communications Systems*, vol. C, no. 24, pp. 333–351, 1994.
- [15] I. C. Paschalidis and J. N. Tsitsiklis, “Congestion-dependent pricing of network services,” *IEEE/ACM Transactions on Networking*, no. 8, pp. 171–184, 1998.
- [16] C. Joe-Wong, S. Ha, and M. Chiang, “Time-dependent broadband pricing: Feasibility and benefits,” in *Proc. of ICDCS*, June 2011.
- [17] S. Zhang, D. Niu, Y. Hu, and F. Liu, “Server selection and topology control for multi-party video conferences,” in *Proc. of ACM NOSSDAV*, Singapore, March 2014.