

Migration Towards Cloud-Assisted Live Media Streaming

Feng Wang, *Member, IEEE*, Jiangchuan Liu, *Senior Member, IEEE*, Minghua Chen, *Member, IEEE*, and Haiyang Wang, *Member, IEEE*

Abstract—Live media streaming has become one of the most popular applications over the Internet. We have witnessed the successful deployment of commercial systems with content delivery network (CDN)- or peer-to-peer-based engines. While each being effective in certain aspects, having an all-round scalable, reliable, responsive, and cost-effective solution remains an illusive goal. Moreover, today's live streaming services have become highly globalized, with subscribers from all over the world. Such a globalization makes user behaviors and demands even more diverse and dynamic, further challenging state-of-the-art system designs. The emergence of *cloud computing*, however, sheds new light into this dilemma. Leveraging the elastic resource provisioning from the cloud, we present Cloud-Assisted Live Media Streaming (CALMS), a generic framework that facilitates a migration to the cloud. CALMS adaptively leases and adjusts cloud server resources in a fine granularity to accommodate temporal and spatial dynamics of demands from live streaming users. We present optimal solutions to deal with cloud servers with diverse capacities and lease prices, as well as the potential latencies in initiating and terminating leases in real-world cloud platforms. Our solution well accommodates location heterogeneity, mitigating the impact from user globalization. It also enables seamless migration for existing streaming systems, e.g., peer-to-peer, and fully explores their potentials. Simulations with data traces from both cloud service providers (Amazon EC2 and SpotCloud) and a live streaming service provider (PPTV) demonstrate that CALMS effectively mitigates the overall system deployment costs and yet provides users with satisfactory streaming latency and rate.

Index Terms—Cloud computing, live media streaming, migration, user/demand globalization.

Manuscript received September 19, 2013; revised March 22, 2014 and August 18, 2014; accepted September 15, 2014; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor L. Ying. Date of publication October 31, 2014; date of current version February 12, 2016. The work of F. Wang was supported by the University of Mississippi under a Start-up Grant. The work of J. Liu was supported by the Canada NSERC under a Discovery Grant. The work of M. Chen was supported by the National Basic Research Program of China under Project No. 2013CB336700 and the University Grants Committee of the Hong Kong Special Administrative Region, China under the Area of Excellence Grant Project No. AoE/E-02/08. The work of H. Wang was supported by the University of Minnesota Duluth under a Start-up Grant. A preliminary version of this paper appeared in the IEEE International Conference on Computer Communications (INFOCOM), Orlando, FL, USA, March 25–30, 2012.

F. Wang is with the Department of Computer and Information Science, The University of Mississippi, University, MS 38677 USA (e-mail: fwang@cs.olemiss.edu).

J. Liu is with the School of Computing Science, Simon Fraser University, Vancouver, BC V5A 1S6, Canada (e-mail: jliu@cs.sfu.ca).

M. Chen is with the Department of Information Engineering, The Chinese University of Hong Kong, Shatin, NT, Hong Kong (e-mail: minghua@ie.cuhk.edu.hk).

H. Wang is with the Department of Computer Science, The University of Minnesota Duluth, Duluth, MN 55812 USA (e-mail: wang4905@d.umn.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNET.2014.2362541

I. INTRODUCTION

IN THE past decade, live media streaming has become one of the most popular applications over the Internet [5], [13]. We have witnessed a number of successful commercial deployments with content delivery network (CDN)- or peer-to-peer-based engines. The former achieves high availability and short startup latencies, but suffers from excessive costs for deploying dedicated servers. This is particularly severe if the user demand fluctuates significantly and the servers have to be overprovisioned for peak loads. The peer-to-peer solution generally incurs lower deployment cost and is more scalable, but the reliability and hence service quality can hardly be guaranteed. There have also been efforts toward synergizing dedicated servers with peer-to-peer [6]. Unfortunately, having an all-round scalable, reliable, responsive, and cost-effective solution remains an illusive goal.

To make it even worse, today's live streaming applications have become highly globalized, with subscribers from all over the world. Such a globalization makes user behaviors and demands even more diverse and dynamic. For illustration, we examine the user demand distribution of PPTV,¹ one of the most popular live media streaming systems in China with multimillion users from a trace analysis [21], [26]. Fig. 1 shows the results of two representative channels (CCTV3 and DragonBall) during one day.² It is easy to see that although PPTV is from China, it has attracted users from all over the world, and the peak time therefore shifts from region to region, depending on the time zone. For example, on the CCTV3 channel, the peak time of North America is around 20:00, while for Asian users, it is around 8:00. During the period of 12:00–20:00, the Asian users have very low demands, while the European users generate most of their demands and the North American users also have moderate demands. Similar observations can also be found from the DragonBall channel, despite that the streaming contents delivered on the two channels are completely different. The impact of such globalized demand turbulence has yet to be addressed in existing systems that largely focus on regional services only.

The emergence of *cloud computing*, however, sheds new light into this dilemma. A cloud platform offers reliable, elastic, and cost-effective resource provisioning, which has been dramatically changing the way of enabling scalable and dynamic network services [4], [10], [19]. There have been studies on demand-driven resource provision [12], [23], [25], [27]; there have been also initial attempts leveraging cloud service to support video-on-demand (VoD) applications, from both industry (e.g.,

¹<http://www.pptv.com>—formerly known as PPLive.

²For ease of comparison, the user demands are normalized by the corresponding maximum demand of each day.

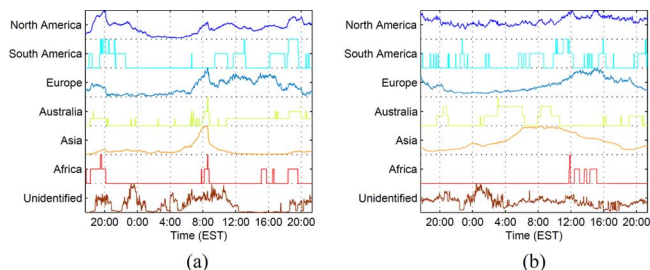


Fig. 1. User demand distributions and variations of a popular live media streaming system (PPTV) on its two typical channels (a) CCTV3 and (b) Drag-onBall during one day.

Netflix [15] and academia [8], [11], [24]. Live media streaming, however, has more stringent playback delay constraints with content being updated in real time. The larger, dynamic, and nonuniform client population further aggravates the problem, calling for new solutions toward a successful migration to the cloud.

In this paper, we present Cloud-Assisted Live Media Streaming (CALMS), a generic framework that facilitates a cost-effective migration to the cloud. CALMS adaptively leases and adjusts cloud servers in a fine granularity to accommodate temporal and spatial dynamics of user demands. We present optimal solutions to deal with cloud servers with diverse capacities and lease prices, as well as the potential latencies in initiating and terminating leases in real-world cloud platforms. Our solution well accommodates location diversity, mitigating the impact from user globalization. It also enables seamless migration for existing streaming systems, e.g., peer-to-peer, and fully explores their potentials. We further develop a set of practical solutions for dynamic adjustment of lease schedules, smart user redirection, and cloud server organization. To understand the performance of CALMS, extensive simulations have been carried out with real data traces from both cloud service providers (Amazon EC2 and SpotCloud) and a live media streaming service provider (PPTV). The results demonstrate that the proposed CALMS effectively mitigates the overall system deployment costs and yet provides users with satisfactory streaming latency and rate.

The remainder of this paper proceeds as follows. Section II presents an overview of the framework. In Section III, we first investigate the basic problem of leasing cloud service and provide an optimal solution, which is then extended by integrating locality awareness and user assistance. The implementation issues and further optimization are discussed in Section IV. We evaluate CALMS by trace-driven simulations in Section V. Section VI further discusses some open issues. Finally, Section VII concludes the paper and discusses potential future directions.

II. CALMS: AN OVERVIEW

Our CALMS intends to provide a generic framework that facilitates the migration of existing live media streaming services to a cloud-assisted solution. Fig. 2 shows an illustration of CALMS, which is divided into two layers, namely *Cloud Layer* and *User Layer*. The Cloud Layer consists of the live media source and dynamically leased cloud servers. Upon receiving a user's subscription request, the Cloud Layer redirects this user to a properly selected cloud server. Such a redirection is transparent to the user, i.e., the Cloud Layer is deemed as a

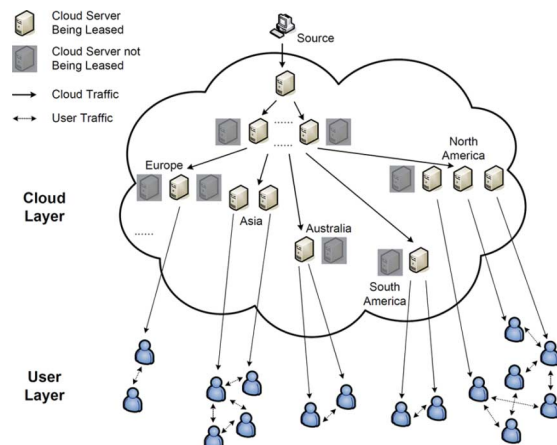


Fig. 2. Overview of CALMS.

single-source server from a user's perspective. Since the user demands change over time, which are also location-dependent, the Cloud Layer accordingly dynamically adjusts the amount and location distribution of the leased servers. Intuitively, it leases more server resources upon demand increase during peak times, and terminates leases upon decrease.

There are, however, a number of critical theoretical and practical issues to be addressed in this generic CALMS framework. First, the cloud servers have diverse capacities and lease prices; the lease duration is not infinitesimally short either, e.g., 1 h for Amazon EC2. As such, when being leased, the server and the pricing cannot be simply terminated at anytime. In addition, for a newly leased server, the configuration and boot up takes time, too, e.g., about 1–10 min for EC2 mainly depending on the used operating system. Though the cloud services are improving, given the hardware, software, and network limits, such latencies can hardly be avoided in the near future. Therefore, CALMS must well predict when to lease a new server to meet the ever-changing demands and when to terminate a server to minimize the lease costs. These problems are further complicated given the global heterogeneous distributions of the cloud servers and that of the user demands.

Note that we do not assume any particular implementation of the User Layer in this study. They can be individual users purely relying on the Cloud Layer, or served by peer-to-peer or a CDN infrastructure, but seeking extra assistance from the Cloud Layer during load surges. In other words, our CALMS framework can smoothly migrate diverse existing live streaming systems. Also, we will explore the potential assistance from user peers in our study by investigating general solutions that well complement existing system designs, taking into account different issues on user dynamics such as user churn, user mobility, and identifying good potential user helpers.

III. FRAMEWORK DESIGN AND SOLUTION

In this section, we discuss the detailed design issues and their solutions for migrating the live media streaming application to the cloud service. Our discussions first start from modeling the basic form of the leasing cloud service problem, and then extend to integrate with locality awareness and user assistance, respectively. We will also present centralized algorithms that yield optimized solutions for addressing these issues, which will further motivate the practical solutions for implementation and optimization in Section IV.

A. Basic Problem: Leasing Cloud Service

Denote the set of cloud servers that can be leased from the cloud service providers as $C = \{c_1, c_2, \dots, c_m\}$. In practice, most cloud providers have a minimum unit time for the duration of leasing a server (e.g., 1 h for Amazon EC2), and when being leased, a cloud server must spend some time in setup and initialization before being ready to use. We use D_m to denote this required minimum unit duration and T_s as the latency for preparing the cloud server. For simplicity, we assume there is always a cloud server directly connecting to the live media source to act as the live media server and use c_0 to denote it. Let R be the rate of the live media streaming. We assume that there is a demand forecast algorithm (such as the algorithms proposed in [16]) to predict the demand in the next period T , where the demand may contain the estimated online population of users, their distributions at different areas or ISPs, and other type of information. At the current stage, we only consider the online population of users and denote it as $\mathbb{P}(t)$ for a given time $t (t \leq T)$. The discussions for utilizing other forecasted demand information will be deferred to later in this section. Define a cloud service lease schedule as $S = \{(x_1, t_1, d_1), (x_2, t_2, d_2), \dots, (x_k, t_k, d_k)\}$ ($t_1 \leq t_2 \leq \dots \leq t_k \leq T$), where a tuple (x_i, t_i, d_i) ($x_i \in C$, $d_i > 0$, and $d_i \bmod D_m = 0$ for $i = 1, 2, \dots, k$) means at time t_i , we start to lease cloud server x_i for the duration d_i . Our problem is thus to find a proper cloud lease schedule S , subject to the following constraints.

1) Service Availability Constraint:

$$\forall (x_i, t_i, d_i) \in S, \text{ if } \exists (x_j, t_j, d_j) \in S \text{ and } x_i = x_j, \\ \text{then } [t_i, t_i + d_i] \cap [t_j, t_j + d_j] = \emptyset.$$

2) Streaming Quality Constraint:

$$\forall T_s \leq t \leq T, \\ \mathbb{U}(c_0) + \sum_{(x_i, t_i, d_i) \in S} (\mathbb{U}(x_i) - R) \cdot I_{[t \in (t_i + T_s, t_i + d_i)]} \geq R \cdot \mathbb{P}(t)$$

where $\mathbb{U}(\cdot)$ is the upload capacity and $I_{[\cdot]}$ is the indicator function. The service availability constraint asks that at any given time, a cloud server can only appear in one schedule. The streaming quality constraint asks that at any given time t , the streaming rate demands of $\mathbb{P}(t)$ users have to be satisfied. Our objective is thus to minimize the lease costs

$$\mathbb{C}_l(c_0) + \sum_{(x_i, t_i, d_i) \in S} \mathbb{C}_l(x_i) \cdot d_i + \mathbb{C}_b(R \cdot \sum_{t \leq T} \mathbb{P}(t))$$

where $\mathbb{C}_l(\cdot)$ and $\mathbb{C}_b(\cdot)$ are the costs for leasing a cloud server and out-cloud bandwidth usage, respectively.³ As the first and last part cannot be reduced, we focus on minimizing the middle part of the total costs, which we denote as $Cost_{\text{lease}}$

$$Cost_{\text{lease}} = \sum_{(x_i, t_i, d_i) \in S} \mathbb{C}_l(x_i) \cdot d_i.$$

³Amazon EC2 had no charges on the traffics into the cloud as well as within the cloud when this research was conducted. Yet, its new policy starts to charge on the traffics among different AWS regions. Our model and solution here can be easily adapted to the new policy by adding an extra cost on the traffic of one streaming rate from outside of the region to inside of the region (except for the region where c_0 is) for each used AWS region.

Algorithm OptimalCloudLeaseSchedule()

```

1: Sort  $C$  by ascendant order of  $\mathbb{C}_l(\cdot)/\mathbb{U}(\cdot)$ ;
2: for  $\forall c \in C$ ,  $leaseTime[c] \leftarrow 0$ ; end for
3: for  $\forall t \leq T$ ,  $Load[t] \leftarrow R \cdot \mathbb{P}(t) - \mathbb{U}(c_0)$ ; end for
4: Set  $Stack$  empty;  $k \leftarrow 0$ ;  $time \leftarrow 0$ ;  $Cost \leftarrow 0$ ;
5:  $renew \leftarrow \text{false}$ ;  $Cost^* \leftarrow \infty$ ; Set  $Stack^*$  empty;
6: do
7:   if  $Cost \geq Cost^*$ , goto 25;
8:   else if  $time + T_s > T$ ,
9:      $Cost^* \leftarrow Cost$ ;  $Stack^* \leftarrow Stack$ ; goto 25;
10:  end if
11:   $k \leftarrow k + 1$ ;
12:  if  $k \leq |C|$  and ( $(!renew \text{ and } leaseTime[c_k] > 0)$  or ( $renew \text{ and } leaseTime[c_k] \neq T_s$ )), continue;
13:  else if  $renew$  and  $k > |C| + 1$  and  $Load[time + T_s] > 0$ ,
14:     $k \leftarrow 0$ ;  $renew \leftarrow \text{false}$ ; continue;
15:  else if ( $!renew$  and  $k > |C|$ ) or ( $renew$  and  $k > |C| + 1$ ), goto 25;
16:  end if
17:  Push  $\{k, time, renew, leaseTime\}$  in  $Stack$ ;
18:  Update  $Load$ ,  $leaseTime$  and  $Cost$ ;
19:  if  $Load[time + T_s] \leq 0$  and ( $!renew$  or  $k = |C| + 1$ ),
20:    Increase  $time$  until  $time + T_s > T$  or  $\exists c \in C$ ,  $leaseTime[c] - \Delta time = T_s$  or  $Load[time + T_s] > 0$ ;
21:     $k \leftarrow 0$ ; Update  $leaseTime$  by  $\Delta time$ ;
22:     $renew \leftarrow \text{if } \exists c \in C, leaseTime[c] = T_s$ ;
23:  end if
24:  continue;
25:  if  $Stack$  is not empty,
26:    Pop  $Stack$ ; Update  $Load$  and  $Cost$ ;
27:  end if
28:  while  $Stack$  is not empty;
29:  Generate optimal cloud lease schedule  $S$  from  $Stack^*$ ;
30:  return  $S$ ;
```

Fig. 3. Algorithm to compute the optimal cloud lease schedule.

This problem is challenging as the cloud servers are scheduled both along the time dimension and at each time instance, along the user demand dimension. By exhaustively searching along both dimensions, the optimal solution can be achieved. However, simply using a naive searching algorithm can be quite inefficient as the solution space increases very quickly with the combination of both dimensions.

To this end, we proposed an enhanced DFS algorithm, which can skip most parts of the solution space and find the optimal solution efficiently. The proposed algorithm is summarized in Fig. 3. To improve the search efficiency, we first sort the cloud servers in C by the ascendant order of the lease cost per unit upload bandwidth (line 1). This allows the servers with cheaper upload bandwidth to be explored first, and near-optimal solutions can then be quickly found. With such solutions, we can further cut other search branches with equal or higher costs (line 7) and greatly reduce the search space for the optimal solution. In addition, we use $leaseTime[c]$ to denote the remaining lease time of server c , and $leaseTime[c] > 0$ means that server c is leased. Since a newly leased server needs T_s for preparation, our algorithm makes decisions in advance of T_s , i.e., at $time$ we lease a new server c so that it can start to provide the service at $time + T_s$. Similarly, for a server c to be renewable at

time, $leaseTime[c]$ must be equal to T_s , so that after we decide to renew it at *time*, it will continue providing the service at *time* + T_s .

We use variable *renew* to distinguish whether the currently considered k th server is renewed or newly leased (lines 11 and 12). When all renewable servers have been considered, we also check if any new server needs to be leased (line 13). If so, we will reset k and *renew* to further explore the branch (line 14). Every time a server c_k is selected (lines 17 and 18), it will be leased for the time of D_m , i.e., $leaseTime[c_k]$ is increased by D_m . $Load[time + T_s]$ to $Load[time + T_s + D_m]$ (or $Load[time + D_m]$, depending on whether c_k is renewed or newly leased) will be reduced by $\mathcal{U}(c_k) - R$. After that, the algorithm checks if any other server needs to be leased or can be renewed (line 19). If not, *time* will be increased (by $\Delta time$) to another position where a server can be renewed or new servers need to be leased. Then, both k and *renew* will be reset accordingly. $leastTime$ will also be updated by subtracting $\Delta time$ for leased servers.

When a search branch is cut off or fully explored, the search will revert to its previous status (lines 25–27), where the previous $\{k, time, renew, leastTime\}$ will be popped up from *Stack* with *Load* and *Cost* being calculated oppositely to line 18. When the search finishes, the optimal cloud lease schedule S will be generated and returned (lines 29 and 30), where a tuple (c_i, t_i, D_m) will be created and added into S for a server c_i newly leased at t_i , and another D_m will be added to the tuple's lease duration if it is renewed afterwards. We then have the following theorem.

Theorem 1: The algorithm proposed in Fig. 3 returns the optimal cloud lease schedule for the basic problem.

Proof: If without cutting branches, our solution will search exhaustively and return the optimal schedule. The proof thus can be done by showing that cutting branches (line 7) do not miss the optimal schedule since *Cost* is a lower bound of those schedules whose search paths contain the current search branch. A more detailed proof can be found in [20]. ■

B. Integrating With Locality Awareness

As the users may be from various locations or time zones over the world, registering to different ISPs and in different areas, a further optimization to the basic problem is thus to integrate with the locality awareness, i.e., to maximize the number of users that connect to the local cloud servers. By achieving this, the delay between users and cloud servers can be effectively reduced, which is often a crucial factor for live media streaming. Another advantage is that such locality awareness can also help to reduce the cross-boundary traffic (e.g., cross-ISP traffic), which is especially important when considering user assistance as discussed in the following. To this end, we use an abstract notation “region” to represent the locality that we care about, which can be interpreted into different meanings in different contexts [e.g., an area with some extent of physical closeness or a group of autonomous systems (ASs) belonging to one ISP]. Let $\mathbb{A} = \{A_1, A_2, \dots, A_n\}$ be the set of regions that the cloud servers and users may be in. For a cloud server c_i , we use $c_i \in A_j$ to denote that it is in the region A_j . In addition, we further extend $\mathbb{P}(t)$ to $\mathbb{P}(A, t)$ to denote the online population of users in region A at time t . The basic problem thus can be rewritten as to find a proper cloud

lease schedule S , subject to the service availability constraint and the following new streaming quality constraint:

$$\begin{aligned} \forall T_s \leq t \leq T, \\ \mathcal{U}(c_0) + \sum_{(x_i, t_i, d_i) \in S} (\mathcal{U}(x_i) - R) \cdot I_{[t \in [t_i + T_s, t_i + d_i]]} \\ \geq R \cdot \sum_{A \in \mathbb{A}} \mathbb{P}(A, t). \end{aligned}$$

Our objective is now to minimize the lease costs

$$Cost_{lease} = \sum_{(x_i, t_i, d_i) \in S} \mathcal{C}_l(x_i) \cdot d_i$$

as well as maximize the locality, which is defined as

$$Locality = \frac{1}{R \cdot \sum_{t \leq T} \sum_{A \in \mathbb{A}} \mathbb{P}(A, t)} \cdot \sum_{t \leq T} \sum_{A \in \mathbb{A}} \min \left(R \cdot \mathbb{P}(A, t), \sum_{(x_i, t_i, d_i) \in S} \mathcal{U}(x_i) \cdot I_{[t \in [t_i + T_s, t_i + d_i], x_i \in A]} \right).$$

These two objectives may contradict with each other, as leasing more servers in each region improves the locality but also raises the lease cost. We adopt the following linear combination form to align them together with different weights:

$$p \cdot \frac{Cost_{lease}}{Cost_{max}} + q \cdot (1 - Locality)$$

where p and q are two parameters that can assign different weights to the two goals. Assuming that the demand estimation is upper-bounded over time duration T , $Cost_{max}$ is thus the minimum lease costs of the case where the demand is constantly set to be the maximum demand within T . As *Locality* is a ratio of the intra-region streaming traffic over the total streaming traffic, $(1 - Locality)$ is thus the ratio of the cross-region traffic over the total traffic, which should also be minimized like $Cost_{lease}$. To make the lease cost part also a ratio ranged between $[0, 1]$, we further divide $Cost_{lease}$ by $Cost_{max}$ and then use parameters p and q to linearly combine the two parts together. This new problem can also be solved by the algorithm proposed in Fig. 3, with some proper but simple modifications. We defer the detailed discussion to Section III-C with the consideration of exploring user resources.

C. Exploring User Resources

It is known that peer-to-peer streaming is highly scalable through exploring user contributed resources. However, in a pure peer-to-peer system, users may join or leave at their own wills, and their available upload bandwidth may vary significantly from time to time and from user to user. Even if the aggregate user contributed upload bandwidth is equal to or greater than the total demand, a pure peer-to-peer design may still suffer from content bottlenecks [14], where users have upload bandwidth but no expected streaming content available for sharing. An extreme is a flashcrowd; that is, during a peak time, many fresh users join the system with a vast amount of ready-to-share upload bandwidth but no content available.

Our CALMS could also benefit from such readily available resources in the User Layer, and yet it can elegantly mitigate the aforementioned problems. To this end, we introduce a group of parameters (α, β) , which we call *user assistance index*. Both of the parameters range from 0 to 1. α determines the ratio of the user contributed upload bandwidth that can be effectively

used to assist the live streaming. β denotes the minimum ratio of cloud servers to be reserved to deal with the content bottleneck (since the cloud servers always have the updated streaming content to share). The full version of our problem thus can be written as to find a proper cloud lease schedule S , subject to the service availability constraint and the following updated streaming quality constraint that considers both locality and user resources:

$$\begin{aligned} \forall T_s \leq t \leq T, \\ \mathcal{U}(c_0) + \sum_{(x_i, t_i, d_i) \in S} (\mathcal{U}(x_i) - R) \cdot I_{[t \in [t_i + T_s, t_i + d_i)]} \\ + \min \left(\alpha \cdot \sum_{A \in \mathbb{A}} \mathbb{B}(A, t), (1 - \beta) \cdot R \cdot \sum_{A \in \mathbb{A}} \mathbb{P}(A, t) \right) \\ \geq R \cdot \sum_{A \in \mathbb{A}} \mathbb{P}(A, t) \end{aligned}$$

where $\mathbb{B}(A, t)$ is the aggregate user upload bandwidth in region A at time t . Our objective is still to minimize

$$p \cdot \frac{Cost_{lease}}{Cost_{max}} + q \cdot (1 - Locality)$$

while with the user assistance taken into account, the locality measure is now calculated as

$$\begin{aligned} Locality = \frac{1}{R \cdot \sum_{t \leq T} \sum_{A \in \mathbb{A}} \mathbb{P}(A, t)} \cdot \sum_{t \leq T} \sum_{A \in \mathbb{A}} \min \left(R \cdot \mathbb{P}(A, t), \right. \\ \left. \sum_{(x_i, t_i, d_i) \in S} \mathcal{U}(x_i) \cdot I_{[t \in [t_i + T_s, t_i + d_i], x_i \in A]} \right) \\ + \min \left(\alpha \cdot \mathbb{B}(A, t), (1 - \beta) \cdot R \cdot \mathbb{P}(A, t) \right) \end{aligned}$$

so as to maximize the traffics from local users/servers, thus reducing the latency and improving the overall performance.

To address this problem, some modifications need to be applied to the algorithm proposed in Fig. 3. First, *Load* now needs to have a second dimension to distinguish the demands from different regions, and when initialized (line 3), it needs to further subtract the user contributed upload bandwidth. In addition, the *Cost* computation (lines 18 and 26) now needs to integrate with the locality, where for a time instance $t \geq time$ (which means that new cloud servers still may be leased for this time instance), we use $R \cdot \mathbb{P}(A, t)$ to approximate the locality in a region A , so that the computed *Cost* still remains a lower bound and the optimality of the cloud lease schedule returned by the algorithm can thus be achieved. We omit more details here, which can be found in [20], due to the space limitation. We then have the following corollary.

Corollary 1: The modified algorithm returns the optimal cloud lease schedule with the consideration of both locality and user resources.

IV. MIGRATION IMPLEMENTATION AND OPTIMIZATION

Section III addresses the major design problems for CALMS. However, in practice, there still remain some issues that need to be considered carefully. First, the user demand and capacity may not be forecasted precisely, where the accuracy often degrades as the predication time gets more forward than the current time. This renders that a statically computed optimal cloud lease schedule for a long period may become less useful due

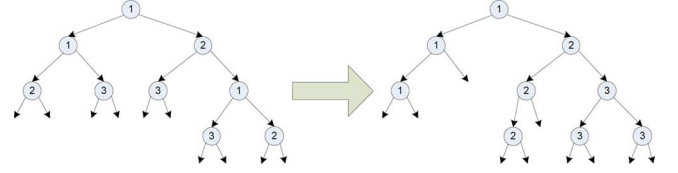


Fig. 4. Cloud Layer organization.

to the prediction error. Another issue is how to organize the leased servers in the Cloud Layer. Since they may be from different regions, a careless organization may introduce unnecessary cross-boundary and cross-region traffics and also degrade the performance. A similar situation also happens to the User Layer, where users need to be carefully organized to enable assistance among each other as well as enforce locality and good quality of service (QoS). In this section, we present our solutions to address these issues.

A. Cloud Layer Organization and Evolution

As the Cloud Layer is the core part of the CALMS framework, its organization is thus very important to the performance of the migrated live media streaming application. As the tree structure is well known for its efficiency for organization, in our implementation, we adopt a tree structure for the Cloud Layer. In particular, we let server c_0 be the root of the tree and in charge of the whole Cloud Layer. The servers in the interior part of the tree relay the live streaming content to other servers to amplify the upload capacity. The remaining servers in outskirts then transmit the live streaming content to the User Layer.

Since the leased cloud servers may be from different regions, naively organizing servers into a tree may still cause poor performance. Fig. 4 shows an example, where each circle denotes a cloud server, and the number inside denotes the region to which the server belongs. In this figure, there are three cloud servers of 2 unit upload capacities leased from each of the regions of 1–3, which are expected to provide similar upload capacities to each region. By a naive organization shown on the left part, the upload capacities provided by cloud servers to each region are 0, 4, and 6, which is greatly different from the expectation; while by handling carefully as shown on the right part, the upload capacities for each region become 3, 3, and 4. In addition, by the naive organization, to arrive at the User Layer, the live media streaming traffics must cross at least one region boundary and may cross as many as three region boundaries; while the carefully handled organization can effectively reduce such cross-boundary traffics.

To deal with these issues, we let server c_0 select and keep tracking a root server for each region from the leased servers. When a cloud server is newly leased, it will first be redirected to the root server in the same region, and the root server will be responsible to help the newly leased server join the subtree rooted at itself. If a region currently has no root server, server c_0 will temporarily take the role until a server from that region is leased. If a root server is stopped due to the decreased demand in its region, the root server will select another server in the same region from its descendants to take its role. If no such server exists, server c_0 will take the role temporarily. In addition, the root server of each region will also help to optimize the server organization within its region, such as allocating the servers with higher upload capacities closer to itself as well as moving a server to a new parent closer to the root server than

the current parent, so that the height of the subtree rooted at the root server can be minimized.

Through our simulation results in Section V, we observe this design, combined with the locality-aware scheduling proposed in Section III, can substantially reduce the cross-boundary traffic as compared to a straightforward approach without locality awareness.

B. User Layer Organization and Evolution

The main task of the User Layer organization is to enforce locality and good QoS as well as enable user assistance. When a user joins the live media streaming session, it will first be redirected to the root server at the same region if there is available upload capacity within that region, otherwise the user will be redirected to the root server of the other region with available upload capacity. The root server then decides where the user should go next. If there is available upload capacity directly from the servers in the region, the user will then be redirected to the server that can provide the upload capacity. If not, the user will be redirected to a server that contains information about users that can provide the upload capacity. This server will then randomly choose some users with enough aggregate upload capacities and send their IPs to the newly joined user. The user will then add these users as neighbors and exchange the live media streaming content with them by a peer-to-peer protocol.

To enforce good QoS and handle the content bottleneck problem that may exist, when a server randomly picks users for the newly joined user, only those with high playback buffer levels will be selected and sent out by the server. In addition, we also let a user with high upload capacity preempt a user that directly downloads from a server but with low upload capacity. When the playback buffer level of a user becomes low, either due to user dynamics or network bandwidth fluctuation, the user will request more neighbors from its server to maintain good live streaming quality.

To track such user information at the Cloud Layer while keeping good scalability, we adopt a hierarchical structure that is naturally provided by the tree organization. We let the last server that the user has been redirected to keep the full information of that user. Such information will then be aggregated and reported to this server's parent. This process will continue on level by level until reaching the root server at that region. The root server of each region will then aggregate and report the user information directly to server c_0 . By this means, when a user is being redirected, the current server that issues the redirection can then choose the next stop for the user based on the collected user information.

To deal with user dynamics, each user will periodically report its updated information to its server. The server will also check whether the user has left if a report has not been received recently. When a server needs to be stopped due to the decreased demands from its region, it will first redirect all its users one by one to server c_0 to redo the join process, then stop and leave the Cloud Layer.

C. Dynamic Lease Scheduling

In practice, the user demand and capacity forecast may be inaccurate, and such inaccuracy may cause a statically computed cloud lease schedule for a long future period to become less

useful or even invalid. To overcome this problem, we use dynamic lease scheduling instead for the migration implementation. In particular, with the collected aggregation information (such as current total user demand and total upload capacity) from each region, server c_0 will dynamically recompute the cloud lease schedule in the following two cases.

- Case 1) If current user demands are greater than the prediction or current user upload capacities are less than the prediction, server c_0 will dynamically recalculate the cloud lease schedule with the updated information and then lease additional cloud servers by the new schedule.
- Case 2) If a cloud server has been leased for the multiple times of D_m (i.e., if necessary, the cloud server can now be stopped without introducing further delay and cost), server c_0 will check if current user demands in that region are less than half of the current upload capacities even with this server stopped. If so, it will stop this server to reduce costs and recompute the cloud lease schedule accordingly.

In addition, due to the content bottleneck issue, sometimes the QoS perceived at users may degrade even though the total upload capacity is still greater than the total user demand. To make the Cloud Layer responsive to such situations, we also collect the playback buffer level of each user when tracking other information. Server c_0 will then check the minimum buffer level of the users who have already started playing the media streaming. If the minimum level is below a threshold, which indicates potential content bottleneck may happen, server c_0 will lease a new server to increase the content availability and recompute a new lease schedule accordingly.

V. PERFORMANCE EVALUATION

We run extensive simulations to evaluate CALMS. The simulations are conducted on the block level and driven by real traces and data sets collected from Amazon EC2 and a popular live media streaming system (PPTV). We first briefly introduce the collected traces and data sets, and then continue to present the methodology as well as the evaluation results.

A. Data Sets and Traces

1) *Amazon EC2 Measurement*: Our measurement on Amazon EC2 is mainly on the bandwidth distribution between cloud servers and users. In particular, we first send DNS requests to the EC2 domains to find the IP addresses of EC2 servers. The DNS server replies with a list of unique IP prefixes that are reserved by Amazon for their EC2 instances. Based on this IP list, we further probe these IP prefixes with ICMP, TCP, and UDP packets from different locations, identifying active instances and then measuring their bandwidth accordingly. Fig. 5 shows the measured bandwidth distributions between different cloud servers and users. Other settings are adopted from previous measurement and evaluation works [3], [10], [17] and the Amazon EC2 official Web site [2]. A more detailed description can be found in [20].

2) *PPTV Traces*: PPTV is one of the largest commercial peer-to-peer live streaming systems to date, attracting over 100 000 online users during peak times, and is also one of the mostly examined systems in academia [5], [7], [26]. Our simulations are based on traces from two of its popular channels, namely

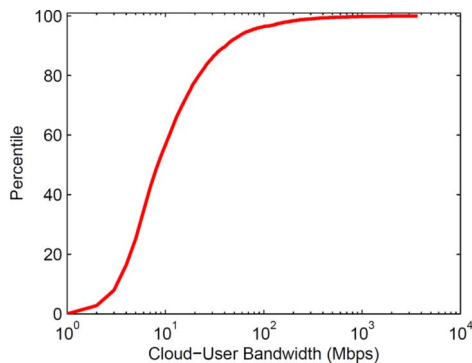


Fig. 5. Measured bandwidth distribution [cumulative distribution function (CDF)] between cloud servers and users.

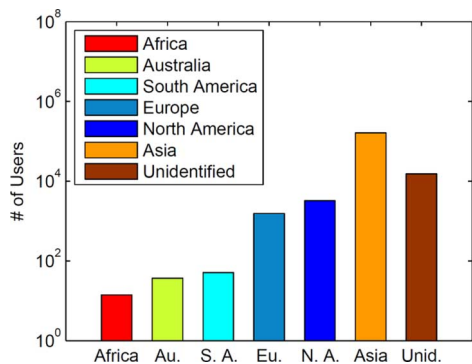


Fig. 6. Region distribution of PPTV users in one day's traces.

CCTV3 and DragonBall. These traces are gathered by an online crawler that continuously collects information from each channel [21]. Fig. 6 shows the region distribution of PPTV users in one day's traces, and Fig. 1 shows the user demands distributions and variations. Other settings about the live media streaming users are adopted from [5] and [26].

B. Methodology

With the data sets and traces from Amazon EC2 and PPTV, we then conduct extensive block-level simulations to evaluate our migration framework. We adopt a typical live media streaming setting as used in PPTV [5] and CoolStreaming [13]. In particular, we set each data block as 1-s video and assume TCP is used for the transmission. The playback buffer is set to the size enough to hold 60 data blocks. The playback will start when at least 30 continuous data blocks are received at the buffer. For comparison, we implement three other approaches: *Max-Central* statically provisions servers all from one region by the maximum user demand in the corresponding trace, which emulates the solution that uses centralized server clusters to provide the live media streaming service. *Max-CDN* also provisions servers by the maximum user demand, but the provisioned servers may be selected from different regions based on the average user demands from each region. This approach emulates the solution that uses CDNs to provide the streaming service. *P2P-Locality* only uses the server c_0 as the streaming server and delivers streaming content by peer-to-peer technology and with locality-awareness, which emulates the solution for peer-to-peer live media streaming.

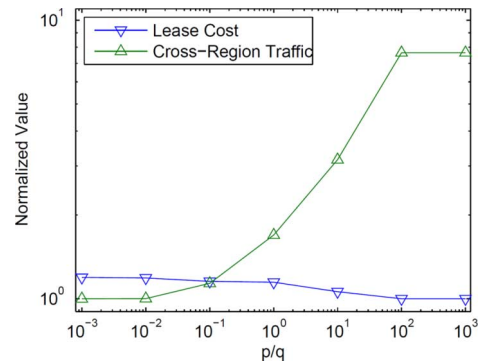


Fig. 7. Impacts of different p/q values on lease cost and cross-region traffic.

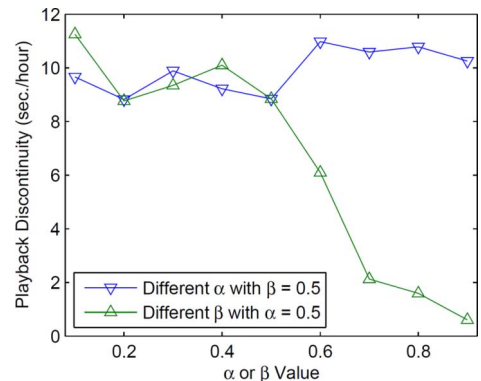


Fig. 8. Impacts of different α or β values on playback discontinuity.

In addition, we use the following four metrics in our simulation: *Lease costs* are the costs for leasing cloud servers, which represent the major concern from the live media streaming service providers. *Cross-region traffic* is the amount of traffic that crosses different regions. This metric shows the locality awareness of an approach and is also an indicator to the general performance as lower cross-region traffic means that users are closer to their servers and the Cloud Layer are well organized. *Playback discontinuity* is the average duration that a user may experience discontinued playback per hour, which is mainly caused by the streaming data packets failing to arrive at a user before its playback deadline. *Startup latency* is the latency taken by a user between its requesting to join the session and receiving enough data to start playback.

C. Impacts of Different Parameter Settings

We first conduct simulations to investigate the impacts of different parameter settings. Fig. 7 demonstrates how the lease cost and cross-region traffic change with different p/q values. For ease of comparison, the results are normalized by the corresponding minimum values. When p/q is small ($\leq 10^{-2}$), the cross-region traffic is minimized while introducing the highest lease costs. On the other hand, when p/q grows large ($\geq 10^2$), it results in the minimum lease cost, but at the expenses of excessive cross-region traffic. Moreover, within the region near 10^{-1} , both the lease cost and cross-region traffic stay relatively low. We thus pick up this value $p/q = 0.1$ as the default for the remaining evaluations.

We next investigate the impacts of different α and β values on the playback discontinuity. The results are shown in Fig. 8. It is easy to see that the impact of β is more significant than that of α ,

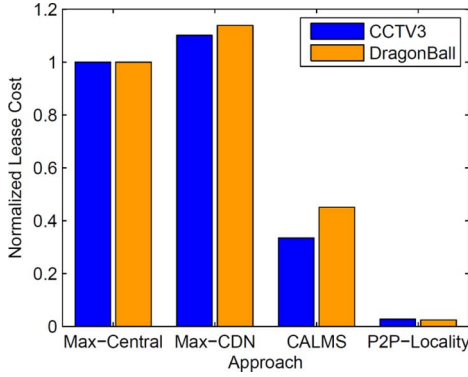


Fig. 9. Lease costs of different approaches.

while the results are different. When α becomes larger, the playback discontinuity slightly increases, which matches the definition of α since more user assistances are involved, resulting in that more content bottlenecks may happen and degrade the playback quality. On the other hand, the playback discontinuity changes inversely with the value of β . This also matches the definition of β as reserving more cloud servers to compensate the content bottlenecks will improve the playback quality. More interestingly, the playback discontinuity will change more dramatically as either of the two parameters changes within the region of $(0.2, 0.7)$. We thus choose $\alpha = 0.2$ and $\beta = 0.7$ as the default setting.

D. Performance Observed at Service Providers

With the default parameter setting, we then conduct simulations to see how CALMS performs with both CCTV3 and DragonBall traces. Fig. 9 shows the lease costs of different approaches. For ease of comparison, the lease costs in each trace are normalized by the corresponding cost of the Max-Central approach. It is not surprising that P2P-Locality has the lowest costs. At the same time, our CALMS also has much lower lease costs, achieving 33.5%–45.1% of the Max-Central approach and 30.5%–39.6% of the Max-CDN approach, respectively. Another observation is that the lease costs in the DragonBall trace are generally higher than those in the CCTV3 trace. This is because the content provided in the DragonBall channel attracts much more Asian user demands than those from the other regions, and the lease prices of Amazon EC2 in Asia are relatively higher, which further verifies the locality-awareness of both the Max-CDN and CALMS approaches. Since Amazon EC2 also charges for the data traffic transferred out of the cloud boundary, we also examine the total costs of different approaches, which is shown in Fig. 10. It is easy to see that even with the data transfer costs, CALMS can still reduce a great amount of the total costs by roughly 28%–30%, which further demonstrates the effectiveness of migrating the live media streaming application to the cloud service.

Fig. 11 shows the cross-region traffics generated by different approaches, which are also normalized by the corresponding cross-region traffic of the Max-Central approach. As the Max-CDN, CALMS, and P2P-Locality approaches are all locality-aware, it is thus not surprising that their cross-region traffics are much lower than that of the Max-Central approach. Yet, an interesting thing is that the cross-region traffics of the three approaches in the DragonBall trace are much lower than those in the CCTV3 trace. This is also because the DragonBall channel

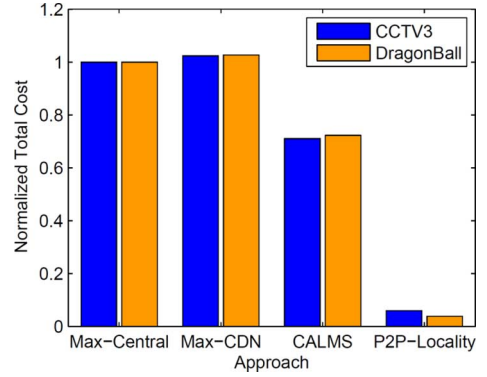


Fig. 10. Total costs of different approaches.

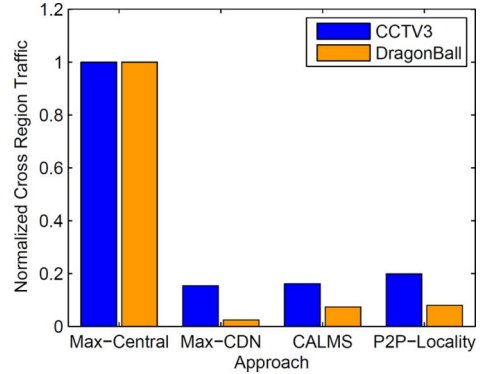


Fig. 11. Cross-region traffics of different approaches.

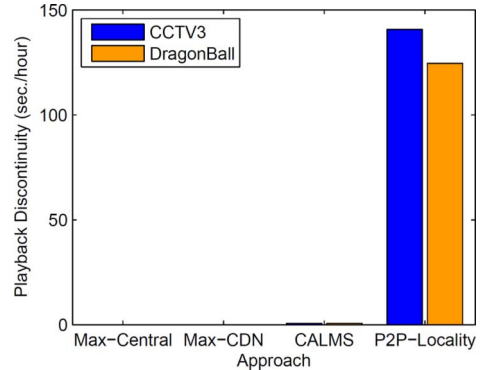


Fig. 12. Playback discontinuity of different approaches.

attracts much more Asian user demands than from the other regions, which allows the server/peer selections to be more dedicated in Asia and results in more intra-region traffics instead of cross-region traffics. Also, by comparing among the three approaches, we can find that by both dynamically provisioning enough cloud servers and exploiting user assistance, CALMS can actually achieve a balance of the locality provided by either mechanism, always avoiding the worst case (having the highest cross-region traffic), and stays close to the best one.

E. QoS Perceived by Users

Section V-D has shown that migrating the live media streaming application to the cloud service can achieve good performance as observed at the service providers. Next, we go on to explore possible benefits that may be brought to users.

We first investigate the playback discontinuity, which is shown in Fig. 12. We can see that CALMS performs similarly

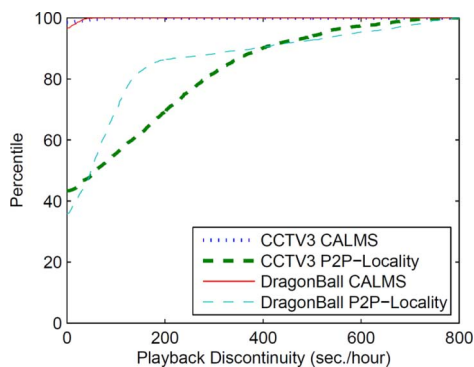


Fig. 13. CDF of playback discontinuity of CALMS and P2P-Locality.

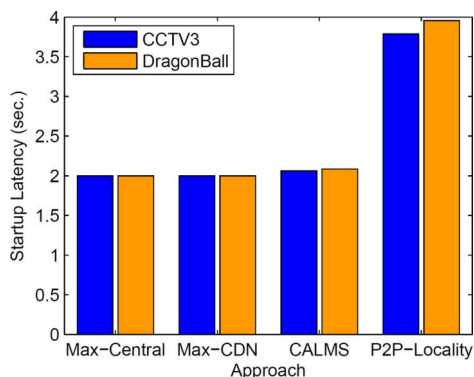


Fig. 14. Startup latency of different approaches.

to both Max- approaches and achieves very low playback discontinuity. On the other hand, due to the peer dynamics and content bottlenecks, P2P-Locality suffers relatively high playback discontinuity with the average as many as over 120 s per hour. To better understand the QoS perceived by each user, we also draw the CDF of the playback discontinuity for both CALMS and P2P-Locality in Fig. 13. It is clear to see that for P2P-Locality, some users may suffer playback discontinuity up to near 800 s per hour, which means that these users cannot watch the live media for more than 20% of the time. On the other hand, by provisioning enough cloud servers to provide both upload contents and capacities, CALMS significantly reduces the worst-case playback discontinuity to around 1 min and affords more than 95% of the users to enjoy zero playback discontinuity.

Besides watching live media fluently, a user also prefers a short waiting time before the live media can start to play. Fig. 14 shows the startup latency of different approaches. Due to the delay of being redirected to other servers and filling up the playback buffers, both Max- approaches have the startup latency of about 2 s. CALMS, by provisioning enough cloud servers with upload contents and capacities, can also achieve a similar startup latency. On the other hand, the startup latency of P2P-Locality is almost of doubled length of the other three approaches due to the long latency for identifying enough peers with both the live media streaming contents and available upload capacities.

VI. FURTHER DISCUSSIONS

As this paper focuses on providing fundamental understandings on migrating the live media streaming to the cloud, we only

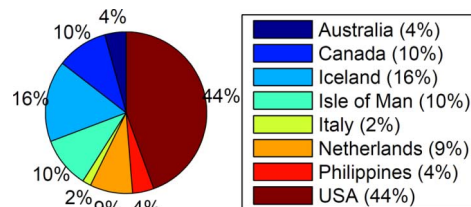


Fig. 15. Location distribution of SpotCloud servers.

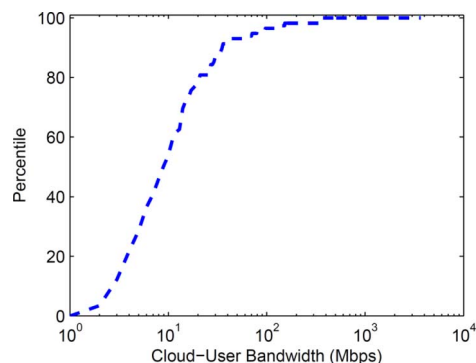


Fig. 16. Measured bandwidth distribution (CDF) between SpotCloud servers and users.

provide basic schemes in Section IV to handle the inaccuracies in the user demand/capacity forecast. It is yet interesting to explore the impact of the forecast accuracy on the performance. Based on our preliminary simulations and analysis, it can be briefly summarized in several folds. First, if the user demand is overestimated, the cross-region traffic and the QoS perceived by users are not affected, while the lease cost increases. If the demand is underestimated, both the lease cost and QoS may drop and the cross-region traffic may increase. On the other hand, if the user capacity is underestimated, both the cross-region traffic and QoS are not affected, while the lease cost may increase. If overestimated, it may degrade the QoS, increase the cross-region traffic, and reduce the lease cost. In CALMS, the capacity forecast inaccuracy is partially overcome by carefully setting α and β , where low α value or high β value can reduce the framework's dependence on the user capacities but may introduce extra lease cost. For the demand forecast error, we believe that to the overall performance, overestimation is better than underestimation. This suggests that one may map the demand forecast problem to the classic ski-rental problem [9] and apply the simple break-even online algorithm. We omit more details here due to the space limitation, which can be found in [20]. Our preliminary results show as long as the decisions of different approaches are made by the same demand forecast algorithm with only overestimation errors, CALMS can still achieve the similar performance as shown in Section V.

Another open issue is that recently increasing attention from both industries and academia have been attracted by a new type of cloud platform featured as allowing the users to contribute/sell their own idle computing resources and build their own cloud services. It is thus interesting to investigate the potential of this type of platform to migrate the live streaming application. To this end, we select a typical example, the Enomaly's SpotCloud [18], [22], to run simulations with the measured information shown in Figs. 15 and 16. The results are shown in Figs. 17–20, where the lease cost and cross-region traffic are

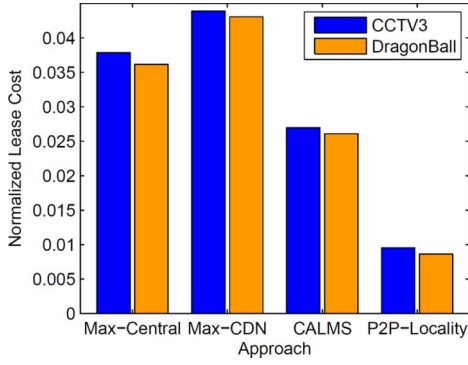


Fig. 17. Lease costs of different approaches by SpotCloud.

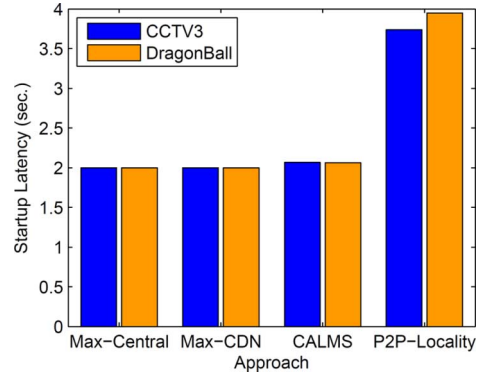


Fig. 20. Startup latency of different approaches by SpotCloud.

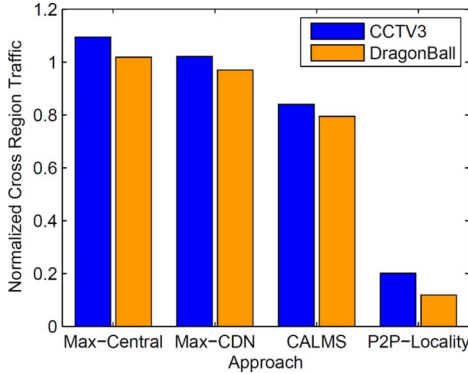


Fig. 18. Cross-region traffics of different approaches by SpotCloud.

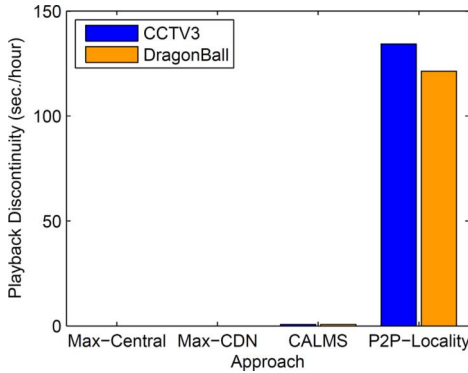


Fig. 19. Playback discontinuity of different approaches by SpotCloud.

normalized by those values of the Max-Central approach on Amazon EC2. Besides the superiority of our CALMS solution, we find two interesting observations compared to the results on Amazon EC2: First, the lease cost of this new type of platform is very low and generally less than 10% of Amazon EC2. This is because this new platform distributes the maintenance cost to the server's contributor/seller, greatly reducing the lease cost of each server. The other observation is that in this new type of platform, although the server contributor/sellers are highly geo-distributed (e.g., eight countries for SpotCloud in Fig. 15), the aggregate server capacities may not be as strong as the data-center clouds (if there is one) in the same region. This explains why the cross-region traffic by SpotCloud is much higher than Amazon EC2. These observations motivate a possible hybrid design to use both types of cloud platforms to provide the live streaming application, where the datacenter clouds are the backbone to enforce enough server capacities and good performance,

with the servers contributed/sold by others from this new type of platform as edge servers to further reduce the lease costs.

VII. CONCLUSION AND FUTURE WORK

This paper presented CALMS, a generic framework that facilitates migrating live media streaming to a cloud platform. Being the very first paper toward this direction, we strived to offer fundamental understandings on the practical feasibility and theoretical constraints in the migration. We first modeled the basic problem of leasing cloud services to support time-varying user demands and developed an optimal algorithm. We then extended our solution to integrate locality awareness and user assistance, alleviating the impact from the service globalization. We further designed a series of practical solutions for both cloud- and user-layer organization and optimization, as well as dynamic lease scheduling that accommodates inaccuracy in demand and capacity forecast. Extensive simulations driven by traces from both cloud service providers (Amazon EC2 and SpotCloud) as well as a live media streaming service provider (PPTV) demonstrated the cost-effectiveness and superior streaming quality of CALMS, even with highly dynamic and globalized demands.

We are currently conducting more simulations to evaluate and improve CALMS with data sets and traces from other cloud providers and live media streaming providers. We expect to develop a prototype for further verification and evaluation. We are also interested in exploring other open issues such as designing better user demand/capacity forecast mechanisms and extending CALMS to other types of cloud platform, e.g., Amazon CloudFront [1], a cloud-based CDN platform.

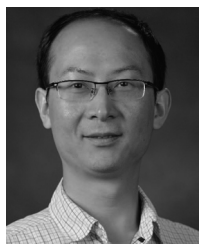
ACKNOWLEDGMENT

The authors thank Dr. X. Hei for providing PPTV traces.

REFERENCES

- [1] Amazon, Seattle, CA, USA, "Amazon CloudFront," [Online]. Available: <http://aws.amazon.com/cloudfront/>
- [2] Amazon, Seattle, CA, USA, "Amazon Elastic Compute Cloud," [Online]. Available: <http://aws.amazon.com/ec2/>
- [3] S. Garfinkel, "An evaluation of Amazon's grid computing services: EC2, S3 and SQS," Harvard University, Cambridge, MA, USA, Tech. Rep., 2008.
- [4] M. Hajjat *et al.*, "Cloudward bound: Planning for beneficial migration of enterprise applications to the cloud," in *Proc. ACM SIGCOMM*, 2010, pp. 243–254.
- [5] X. Hei, C. Liang, J. Liang, Y. Liu, and K. Ross, "A measurement study of a large-scale P2P IPTV system," *IEEE Trans. Multimedia*, vol. 9, no. 8, pp. 1672–1687, Dec. 2007.

- [6] C. Huang, J. Li, and K. W. Ross, "Peer-assisted VoD: Making Internet video distribution cheap," in *Proc. USENIX IPTPS*, 2007, pp. 1–6.
- [7] Y. Huang, T. Z. J. Fu, D.-M. Chiu, J. C. S. Lui, and C. Huang, "Challenges, design and analysis of a large-scale P2P-VoD system," in *Proc. ACM SIGCOMM*, 2008, pp. 375–388.
- [8] Z. Huang, C. Mei, L. E. Li, and T. Woo, "Cloudstream: Delivering high-quality streaming videos through a cloud-based SVC proxy," in *Proc. IEEE INFOCOM Mini-Conf.*, 2011, pp. 201–205.
- [9] A. Karlin, M. Manasse, L. Rudolph, and D. Sleator, "Competitive snoopy caching," *Algorithmica*, vol. 3, no. 1, pp. 79–119, 1988.
- [10] A. Li, X. Yang, S. Kandula, and M. Zhang, "CloudCmp: Comparing public cloud providers," in *Proc. ACM IMC*, 2010, pp. 1–14.
- [11] H. Li, X. Cheng, and J. Liu, "Understanding video sharing propagation in social networks: Measurement and analysis," *Trans. Multimedia Comput., Commun. Appl.*, vol. 10, no. 4, p. 33, 2014.
- [12] M. Lin, A. Wierman, L. L. H. Andrew, and E. Thereska, "Dynamic right-sizing for power-proportional data centers," *IEEE/ACM Trans. Netw.*, vol. 21, no. 5, pp. 1378–1391, Oct. 2013.
- [13] J. Liu, S. G. Rao, B. Li, and H. Zhang, "Opportunities and challenges of peer-to-peer Internet video broadcast," *Proc. IEEE*, vol. 96, no. 1, pp. 11–24, Jan. 2008.
- [14] N. Magharei and R. Rejaie, "PRIME: Peer-to-peer Receiver-driven MESH-based streaming," in *Proc. IEEE INFOCOM*, 2007, pp. 1415–1423.
- [15] Netflix, Los Gatos, CA, USA, "NetFlix," [Online]. Available: <http://www.netflix.com>
- [16] D. Niu, Z. Liu, B. Li, and S. Zhao, "Demand forecast and performance prediction in Peer-Assisted On-Demand streaming systems," in *Proc. IEEE INFOCOM Mini-Conf.*, 2011, pp. 421–425.
- [17] R. Shea, F. Wang, H. Wang, and J. Liu, "A deep investigation into network performance in virtual machine based cloud environment," in *Proc. IEEE INFOCOM*, 2014, pp. 1285–1293.
- [18] SpotCloud, Bethesda, MD, USA, "SpotCloud," [Online]. Available: <http://www.spotcloud.com/>
- [19] K. Sripanidkulchai, S. Sahu, Y. Ruan, A. Shaikh, and C. Dorai, "Are clouds ready for large distributed applications?," *Oper. Syst. Rev.*, vol. 44, no. 2, pp. 18–23, Apr. 2010.
- [20] F. Wang, J. Liu, M. Chen, and H. Wang, "Cloud-assisted live media streaming," University of Mississippi, University, MS, USA, Tech. Rep., 2014.
- [21] F. Wang, Y. Xiong, and J. Liu, "mTreebone: a collaborative tree-mesh overlay network for multicast video streaming," *IEEE Trans. Parallel Distrib. Syst.*, vol. 21, no. 3, pp. 379–392, Mar. 2010.
- [22] H. Wang, F. Wang, J. Liu, and J. Groen, "Measurement and utilization of customer-provided resources for cloud computing," in *Proc. IEEE INFOCOM*, 2012, pp. 442–450.
- [23] C. Wu, B. Li, and S. Zhao, "Dynamic server provisioning in multi-channel P2P live streaming," *IEEE/ACM Trans. Netw.*, vol. 19, no. 5, pp. 1317–1330, Oct. 2011.
- [24] Y. Wu, C. Wu, B. Li, X. Qiu, and F. C. Lau, "CloudMedia: When cloud on demand meets video on demand," in *Proc. IEEE ICDCS*, 2011, pp. 268–277.
- [25] H. Xu and B. Li, "Joint request mapping and response routing for geo-distributed cloud services," in *Proc. IEEE INFOCOM*, 2013, pp. 854–862.
- [26] K. Xu, H. Li, J. Liu, W. Zhu, and W. Wang, "PPVA: A universal and transparent peer-to-peer accelerator for interactive online video sharing," in *Proc. IWQoS*, 2010, pp. 1–9.
- [27] Q. Zhang, M. F. Zhani, R. Boutaba, and J. L. Hellerstein, "HARMONY: Dynamic heterogeneity-aware resource provisioning in the cloud," in *Proc. IEEE ICDCS*, 2013, pp. 510–519.

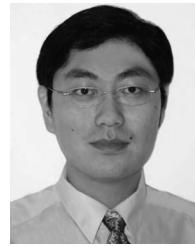


Feng Wang (S'07–M'13) received the Bachelor's and Master's degrees in computer science and technology from Tsinghua University, Beijing, China, in 2002 and 2005, respectively, and the Ph.D. degree in computing science from Simon Fraser University, Burnaby, BC, Canada, in 2012.

He is currently an Assistant Professor with the Department of Computer and Information Science, The University of Mississippi, University, MS, USA. His research interests include cloud computing, peer-to-peer networks, socialized content sharing, wireless

mesh/sensor networks, and cyber-physical systems.

Dr. Wang is a TPC Co-Chair in IEEE ICC 2014 for the Symposium on Wireless Networking and Multimedia.



Jianguan Liu (S'01–M'03–SM'08) received the B.Eng. degree (*cum laude*) from Tsinghua University, Beijing, China, in 1999, and the Ph.D. degree from The Hong Kong University of Science and Technology, Hong Kong, in 2003, both in computer science.

He is a Full Professor with the School of Computing Science, Simon Fraser University, Vancouver, BC, Canada, and an EMC-Endowed Visiting Chair Professor with Tsinghua University. He was an Assistant Professor with the Department of Computer Science and Engineering, The Chinese University of Hong Kong, Hong Kong, from 2003 to 2004. His research interests include multimedia systems and networks, cloud computing, social networking, wireless ad hoc and sensor networks, and peer-to-peer and overlay networks.

Prof. Liu serves on the editorial boards of the IEEE TRANSACTIONS ON MULTIMEDIA, IEEE COMMUNICATIONS SURVEYS AND TUTORIALS, IEEE ACCESS, and IEEE INTERNET OF THINGS JOURNAL. He is a TPC Co-Chair of IEEE/ACM IWQoS 2014 and an Area Chair of ACM Multimedia 2014. He is a co-recipient of the ACM TOMCCAP Nicolas D. Georganas Best Paper Award 2013, ACM Multimedia Best Paper Award 2012, IEEE GLOBECOM 2011 Best Paper Award, and IEEE Communications Society Best Paper Award on Multimedia Communications 2009.

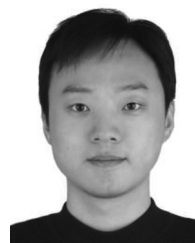


Minghua Chen (S'04–M'07) received the B.Eng. and M.S. degrees in electronic engineering from Tsinghua University, Beijing, China, in 1999 and 2001, respectively, and the Ph.D. degree in electrical engineering and computer sciences from the University of California (UC), Berkeley, CA, USA, in 2006.

He spent one year visiting Microsoft Research, Redmond, WA, USA, as a Postdoctoral Researcher. He joined the Department of Information Engineering, The Chinese University of Hong Kong,

Hong Kong, in 2007, where he currently is an Associate Professor. He is also an Adjunct Associate Professor with the Peking University Shenzhen Graduate School, Shenzhen, China, in 2011–2014. His recent research interests include energy systems (e.g., microgrids and energy-efficient data centers), distributed optimization, multimedia networking, wireless networking, network coding, and secure network communications.

Dr. Chen is currently an Associate Editor of the IEEE/ACM TRANSACTIONS ON NETWORKING. He received the Eli Jury Award from UC Berkeley (presented to a graduate student or recent alumnus for outstanding achievement in the area of systems, communications, control, or signal processing) in 2007 and The Chinese University of Hong Kong Young Researcher Award in 2013. He also received several Best Paper awards, including the IEEE ICME Best Paper Award in 2009, the IEEE TRANSACTIONS ON MULTIMEDIA Prize Paper Award in 2009, and the ACM Multimedia Best Paper Award in 2012.



Haiyang Wang (S'08–M'13) received the Ph.D. degree in computing science from Simon Fraser University, Burnaby, BC, Canada, in 2013.

He is currently an Assistant Professor with the Department of Computer Science, University of Minnesota Duluth, Duluth, MN, USA. His research interests include cloud computing, big data, socialized content sharing, multimedia communications, and peer-to-peer networks.