

Toward Cloud-Based Distributed Interactive Applications: Measurement, Modeling, and Analysis

Haiyang Wang, Tong Li¹, Ryan Shea, Xiaoqiang Ma, Feng Wang, *Member, IEEE*,
Jiangchuan Liu, *Fellow, IEEE*, and Ke Xu, *Senior Member, IEEE, ACM*

Abstract—With the prevalence of broadband network and wireless mobile network accesses, distributed interactive applications (DIAs) such as online gaming have attracted a vast number of users over the Internet. The deployment of these systems, however, comes with peculiar hardware/software requirements on the user consoles. Recently, such industrial pioneers as Gaikai, Onlive, and Ciinow have offered a new generation of cloud-based DIAs (CDIAs), which shifts the necessary computing loads to cloud platforms and largely relieves the pressure on individual user’s consoles. In this paper, we aim to understand the existing CDIA framework and highlight its design challenges. Our measurement reveals the inside structures as well as the operations of real CDIA systems and identifies the critical role of cloud proxies. While its design makes effective use of cloud resources to mitigate client’s workloads, it may also significantly increase the interaction latency among clients if not carefully handled. Besides the extra network latency caused by the cloud proxy involvement, we find that computation-intensive tasks (e.g., game video encoding) and bandwidth-intensive tasks (e.g., streaming the game screens to clients) together create a severe bottleneck in CDIA. Our experiment indicates that when the cloud proxies are virtual machines (VMs) in the cloud, the computation-intensive and bandwidth-intensive tasks may seriously interfere with each other. We accordingly capture this feature in our model and present an interference-aware solution. This solution not

only smartly allocates workloads but also dynamically assigns capacities across VMs based on their arrival/departure patterns.

Index Terms—Cloud-based distributed interactive application, interaction latency, task interference.

I. INTRODUCTION

WITH unmatched system flexibility, distributed interactive applications (DIAs) have become increasingly popular in recent years. By providing diverse interactions among users, such applications as massive multiplayer online gaming, live messaging, and shared whiteboard have attracted a vast number of users over the Internet. Take online gaming as an example. It is reported in [1] that nowadays each US household on average owns at least one dedicated game console or PC for game playing, where 62% of them have played interactive games with others. The global markets of these DIA systems will also expand from 58.7 billion in 2011 to 83 billion in 2016, growing at a 7.2% compound annual rate [2]. Yet, to support superior interactions, DIAs often have peculiar demands on user consoles. The specialized consoles with high-performance hardware unavoidably increase user cost and greatly limit the penetration of DIAs to ubiquitous end users.

To bring true *play-as-you-go* into reality, industrial pioneers such as Gaikai [3], Onlive [4] and Ciinow [5] have suggested a new generation of DIAs based on cloud computing platforms. This *cloud-based distributed interactive application* (CDIA) effectively shifts the hardware/software requirements as well as the necessary computing loads to *cloud proxies*, and thus has attracted an increasing amount of attention from both service providers and end users. In particular, Sony Computer Entertainment (SCE) acquired Gaikai for 380 million USD on July 2, 2012, putting Gaikai as a key function in its next generation game console, PlayStation 4 [6]. Its competitor, Microsoft, also announced that the Gaikai-like CDIA functions will also play a major role in their new game console Xbox One [7]. Moreover, such industry leaders as AMD and Nvidia are also entering the market of CDIA services [8]. AMD’s investment in CiiNow gives it a means of competing with rival Nvidia’s GeForce Grid cloud gaming platform. In addition to cloud gaming, CDIAs can also take the form of the live notebook environments (e.g., Jupyter [9]) or a dashboard with live interactive graphical widgets (e.g., web applications based on R’s shiny package [10]). These CDIAs also shift the software requirements to the cloud, and offload the compiling and parsing workloads to high-performance machines, i.e., cloud proxies.

Manuscript received July 5, 2016; revised March 9, 2017 and July 16, 2017; accepted September 22, 2017; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor Y. Ganjali. Date of publication November 14, 2017; date of current version February 14, 2018. This work was supported in part by the National Natural Foundation of China under Grant 61472212, in part by the National Science and Technology Major Project of China under Grant 2015ZX03003004, in part by the National High Technology Research and Development Program of China 863 Program, in part by the EU Marie Curie Actions CROWN under Grant FP7-PEOPLE-2013-IRSES-610524, in part by the NSERC Discovery Grant, in part by the E.W.R. Steacie Memorial Fellowship, and in part by the Industrial Canada Technology Demonstration Program Grant. The work of H. Wang was supported by the Chancellors Small Grant and Grant-in-Aid Programs from the University of Minnesota. The work of F. Wang was supported by a Start-up Grant from the University of Mississippi. (*Corresponding authors: Ke Xu; Tong Li.*)

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

T. Li is with Huawei Technologies Co., Ltd., Shenzhen 518129, China (e-mail: litong12@tsinghua.org.cn).

R. Shea and J. Liu are with the Department of Computer Science, Simon Fraser University, Burnaby, BC V5A 1S6, Canada (e-mail: rws1@cs.sfu.ca; jliu@cs.sfu.ca).

X. Ma is with the School of Electronic Information and Communications, Huazhong University of Science and Technology, Wuhan 430073, China (e-mail: mxqhust@gmail.com).

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

K. Xu is with the Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China (e-mail: xuke@tsinghua.edu.cn).

This paper has supplementary downloadable material available at <http://ieeexplore.ieee.org>, provided by the authors.

Digital Object Identifier 10.1109/TNET.2017.2765246

Nowadays, CDIA remains in its infancy with plenty of unknown issues. In this paper, we take a first step towards understanding the CDIA framework and highlight its design challenges. Our measurement reveals the inside structure as well as the operations of real CDIA systems and identifies the critical role of cloud proxies. While this design makes effective use of cloud resources to mitigate client's workloads, it may also significantly increase the interaction latency among clients if not carefully handled. In both DIA and CDIA systems, the interaction latency is the most essential problem even when there is no limitation on the availability of server capacities [11] [12]. The increasing latency will unavoidably reduce the applicability of CDIA systems.

In detail, our analysis reveals that the deployment of cloud proxies adds extra communication hops between clients. To make the matter worse, the processing latency at the cloud proxy is also surprisingly high. While the use of the high-performance cloud platforms is expected to be highly efficient, we find that the computation-intensive tasks (e.g., game video encoding) and the bandwidth-intensive tasks (e.g., streaming game screens to clients) together create a severe bottleneck in CDIAs. Our experiment indicates that when the cloud proxies are virtual machines (VMs) in the cloud, the computation-intensive and bandwidth-intensive tasks may seriously interfere with each other if not handled carefully. An increase of traffic load will greatly slow down the CPU benchmark of cloud VMs. In the case of CDIA, when the cloud proxies are used to stream game screens to users, the computation-intensive operations, such as game processing and message forwarding, will also be invoked and prolong the interaction latency. The large number of CDIA users further aggravates this issue with mutual-interference, leading to poor user experiences. For example, *Diablo 2* in USEast realm has over 5 million users across 20 servers. The maximum acceptable latency of this Role Playing Games (RPG) is from 100 to 500 ms [13]. However, when we deploy these systems on CDIA, their latency can easily exceed 600 ms even when we assume that the network latency is as small as zero.

Such interference however does not exist in conventional physical machines or to a much lower degree. Therefore, the existing load assignment solutions in the DIA system are mainly focused on the optimization of stand-alone workloads, without considering their interference in the VM environment. To address this problem, we provide a model and consider an interference-aware solution that not only smartly allocates workloads but also dynamically assigns capacities across different VMs. Moreover, this model is also enhanced to capture user's arrival/departure patterns. The evaluation indicates that the proposed enhancement can successfully handle user dynamic while providing a bounded interaction latency.

The rest of this paper is organized as follows. Based on the measurement of Section II, we examine the latency issues in CDIA in Section III. Section IV reveals the task interference on cloud proxies and Section V provides a model to minimize the interaction latency. Our model is then extensively examined in Section VI. Section VII provides the solution to consider user's arrival/departure dynamic. Section VIII presents the related work and Section IX concludes the paper.

II. CLOUD-BASED DIA: BACKGROUND AND FRAMEWORK

From the industrial perspective, CDIAs can bring immense benefits by expanding user bases to the vast number of less-powerful devices that support thin clients only, particularly smartphones and tablets. As an example, the recommended system configuration for *Battlefield 3*, a highly popular first-person shooter game, is a quad-core CPU, 4 GB RAM, 20 GB storage space, and a graphics card with at least 1GB RAM (e.g., NVIDIA GEFORCE GTX 560 or ATI RADEON 6950), which alone costs more than \$500 USD. The newest tablets (e.g., Apple's iPad with Retina display and Google's Nexus 10) cannot even meet the minimum system requirements that need a dual-core CPU over 2.4 GHz, 2 GB RAM, and a graphics card with 512 MB RAM, not to mention smartphones of which the hardware is limited by their smaller size and thermal control. Furthermore, mobile terminals have different hardware/software architectures from PCs, e.g., ARM rather than x86 for CPU, lower memory frequency and bandwidth, power limitations, and distinct operating systems. Therefore, the traditional console game model is unfeasible for these devices, which in turn become targets for such CDIA systems as Gaikai and Onlive. Different from existing DIAs, CDIAs utilize the powerful and elastic service capacity offered by cloud computing to mitigate the hardware/software requirements on user consoles. In particular, Gaikai and Onlive deploy the game clients/consoles on cloud platforms and only stream the game screen/interactions to end users.

To understand how CDIAs work in detail, we take Gaikai as a case study. Since 2011, it has emerged as one of the most popular cloud-based online gaming systems with over 100 million subscribers. It not only provides free PC game demos but also powers high quality gaming experiences onto smartphones, tablets and Internet TVs [14].

Our experiments have conducted traffic measurement and analysis from the edge of four networks, which are located in four different countries (United States, Canada, China and Japan) in two distinct continents. We used traffic analysis, shared library and system call interception techniques to analyze various aspects of the Gaikai protocol. We also applies a packet level analysis to understand the communications between our clients and cloud proxies. In particular, we monitor Gaikai's online gaming service with clients from different network locations and capture their traffic from/to the Gaikai servers. We first examine the details in the control messages, look into the packet payloads and identify the information, such as the domain names, of Gaikai's basic components. After that, we carefully analyze the data traffic to further understand the functions of these components. To avoid possible bias, we apply classic analysis tools (e.g., ISP lookup) that are widely adopted in the existing reverse engineering studies [15]. Note that some online games are designed to utilize one of the user clients as the server to enable interactions. We therefore take advantage of this feature and capture the server-side traffic to better understand the related protocols.

Through analyzing the captured traffic, we illustrate Gaikai's basic framework/protocol in Figure 1. We can see that there are two major components on the server side of

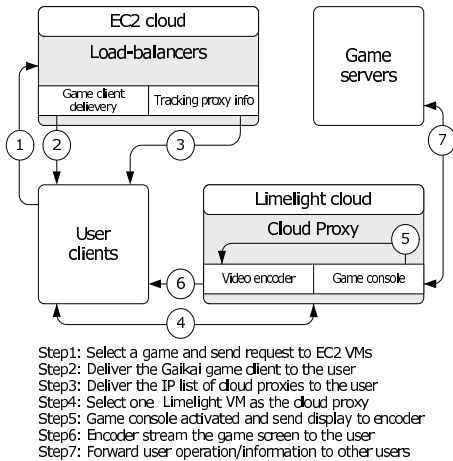


Fig. 1. Basic Framework of Gaikai.

Gaikai (marked as grey boxes in the figure). The first part is Amazon EC2-based [16] load-balancers,¹ and the second part is the Limelight-based game proxy servers [17]. Both Amazon and Limelight are leading cloud service providers with Xen virtualization [18]. Gaikai applies both platforms to accomplish different functionalities and utilizes their widely geo-distributed instances to push these functions closer to the users.

When a user selects a game on Gaikai (*Step 1* in Figure 1), an EC2 virtual machine (VM) will first deliver the Gaikai game client to the user (in *Step 2*). After that, it forwards the IP addresses of the available Limelight game proxies to the users (in *Step 3*). The user will then use one game proxy to run the game (in *Step 4*). To ensure smooth game playing, this selected game proxy uses a *packet train measurement* [19] to estimate the available bandwidth to the users. This is identified by our packet-level analysis, which shows that the game proxy sends back-to-back packets with empty payload to test the available bandwidth. Note that the game proxy starts the game only when the available bandwidth can well-support an FPS (frames per second) around 60 for video streaming. After that, the game proxy starts to run the game and the game screen will be streamed back to the user via UDP (in *Step 5* and *Step 6*). For multiplayer online games, these game proxies will also forward user operations to the game servers (mostly deployed by game developers) and send the related information/reactions back to users (in *Step 7*). It is easy to see that such a CDIA system can remarkably relieve the hardware/software requirements on the user side, given that now the games are running on the cloud platforms. This change enables users to play hard-core games over much less powerful devices, e.g., over smartphone, tablets, or even digital TVs, as long as they are multimedia- and network-ready.

We have also measured other CDIA platforms, and have found that Gaikai’s framework is representative, which is not surprising given it as a very natural extension to the conventional DIA with cloud assistance.

¹Based on our measurement, they also have other functions beside load-balancing. We call them *load-balancers* because Gaikai marks them with “LB” in their domain names.

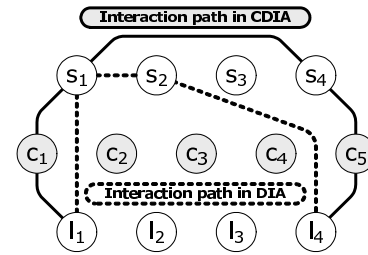


Fig. 2. Path of client interaction.

III. INTERACTION LATENCY OF CDIA: ISSUES AND CHALLENGES

The CDIA framework offers great opportunities for both users and service providers. Similar to the conventional DIA systems, its service performance is highly depending on users’ *interaction latency* [12]. In particular, the CDIA systems, such as the online cloud gaming applications, need to collect user actions, transmit them to the cloud proxy, process the action, render the results, encode/compress the resulting changes to the game-world, and stream the video (game scenes) back to the player. To ensure interactivity, all of these serial operations must happen in the order of milliseconds. The interaction latency is thus essential even when there are no limitations on the availability of server capacities [11].

Figure 2 illustrates the interaction paths in both DIA and CDIA frameworks, where L is the set of clients, S is the set of servers, and C is the set of cloud-based proxies. To better compare the performance of DIA and CDIA design, we will first focus on the latency between clients and servers (e.g., between l_1 and s_1). It is easy to see that such a latency consist of two parts in the CDIA systems: First, the network latency between clients and servers; Second, the processing latency on cloud proxies. We will provide a step-by-step discussion to understand these two parts in the following subsections.

A. Network Latency Analysis

To clarify the extra network latency in CDIA design, we carry out a real-word experiment from Planet-lab. We use a server in our campus to emulate the game server in CDIA.² We select 588 Planet-lab nodes (the maximum number of nodes that we can access) to run as CDIA clients and emulate the CDIA framework by using the server and these clients to connect Gaikai’s cloud proxies. We have found 28 cloud proxies during the measurement of Gaikai,³ and therefore use the IP addresses of these proxies in this experiment. We first measure the RTTs between 588 Planet-lab clients and 28 Gaikai cloud proxies and then the RTTs between the server and these cloud proxies. The sum of these two latencies can be used to calculate the client-server RTTs in this CDIA system. To provide a fair comparison, we also measure the direct RTTs between the server and the Planet-lab clients as the baseline (the case of conventional DIA).

²We have observed similar results over 50 servers that are located in different places.

³The total number of Gaikai’s cloud proxy is unknown to the general public. These sampled cloud proxies are used to estimate the network latency in such a system.

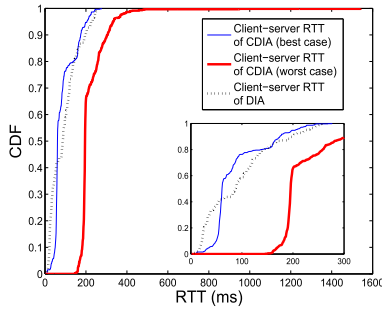


Fig. 3. Time cost between user and server (DIA v.s. CDIA).

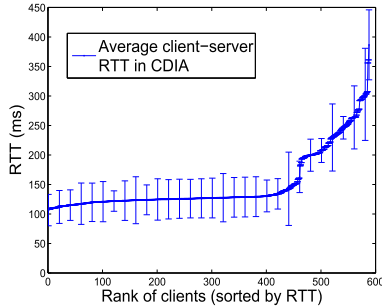


Fig. 4. Average user-server latency in CDIA.

Figure 3 compares the client-server RTT in both DIA and CDIA. We can see that most (over 80%) users in DIA have quite low client-server latency (less than 60 ms), while the average latency is much worse if we put them into CDIA, as shown in Figure 4. The worst case in Figure 4 shows 90% users have an interaction latency over 200 ms. This is hardly acceptable for smooth interaction because the processing latencies are not yet considered in this experiment. It is known that adding extra nodes in any overlay network may unnecessarily lead to longer path length given that triangle inequality does not hold in the Internet [20]. Hence, there is indeed space to reduce the latency beyond naive proxy deployment.⁴

B. Processing Latency Analysis

So far, we have only considered the network latency in the CDIA framework. It is easy to see that the cloud proxies will also bring extra processing latency to the interaction. We now closely examine this latency and identify its impacts.

To focus exactly on the interaction between clients and cloud proxies, we select a single player game where a player (client) does not need to communicate with the game server and other players. The player simply sends the operations to the cloud proxy and the proxy then streams the responding game screen back to the player. Since the RTT between the player and the cloud proxy can be directly measured, we only need to obtain the *response time*,⁵ which, after subtracting

⁴The 588 PlanetLab nodes are applied in both DIA and CDIA experiments to provide fair comparison. Since some real-world interactive applications, such as Diablo 2, may divide their users into realms, we also investigate the case using a subset of PlanetLab nodes from one realm (e.g., USEast). The results remain consistent to Figure 3.

⁵The response time is the latency that the player waits until the result of her/his operations is returned. For example, if the player clicks the button “option” at time t_i and the option menu displays at time t_j , the response time can be calculated as $t_j - t_i$.

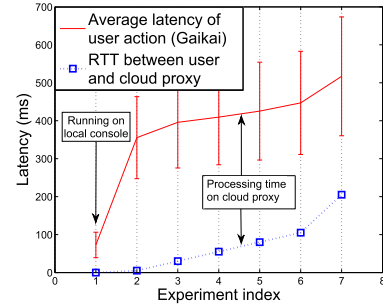


Fig. 5. Average latency of user’s action.

the RTT, gives the processing latency at the cloud proxy. The detail of this experiment is as follows.

We first select an action button in the game *The Witcher 2: Assassins of Kings*; in particular, the “map” button. We click this button and start to record the game screen at 100 FPS (frames per second). This sampling rate already exceeds the normal game play which is around 60 to 70 FPS. We then check the video file frame-by-frame until we find the frame where the map is displayed. We run this experiment for multiple times under different RTTs. These RTTs are controlled by the traffic shaping tool *TC* [21]. To better understand the processing overhead at the cloud proxy, we also record the response time on a local game console. We use the same game on both Gaikai and the local console to provide a fair comparison.

As we can see from Figure 5, the local console general needs 80 ms to open the map for the players with very small standard deviation. Note that the RTT is zero in this case because the game is locally rendered. When we run this game remotely on Gaikai, the response time elevates to more than 300 ms. The overhead (in terms of the processing latency) on the Gaikai proxy is thus approximately 220 ms.⁶ When we consider the interactive latency between different users, there will be, unfortunately, two proxies in their interactive path. This means the interactive latency can easily exceed 600 ms. It is worth noting that the studies on traditional gaming systems have found that different styles of games have different thresholds for maximum tolerable delays [13], [22]–[24]. These thresholds can be further adjusted according to the user feedback and operating experience, this is because that indeed the slower the game is, the less delay influences the user rating [25]. However, latency as high as 600 ms will not be acceptable for most online gaming applications as shown in Table I.

IV. INTERFERENCES BETWEEN COMPUTATION-INTENSIVE AND BANDWIDTH-INTENSIVE TASKS

It is surprising to see that the cloud proxies can introduce such a high processing latency in CDIA. This is unlikely caused by video encoding only since many CDIA service providers have claimed that their video encoding latency is

⁶To avoid measurement bias, we also test actions that make different changes to the in game world, for example, small character movements. The results remain consistent with Figure 5.

TABLE I
DELAY TOLERANCE IN TRADITIONAL GAMING

Example Game Type	Perspective	Delay Threshold
First Person Shooter (FPS)	First Person	100 ms
Role Playing Game (RPG)	Third-Person	500 ms
Real Time Strategy (RTS)	Omnipresent	1000 ms

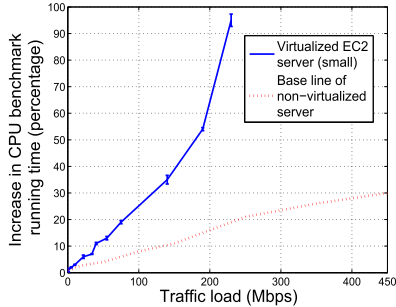


Fig. 6. Increase of CPU benchmark running time (on EC2 small instance).

indeed very small within 10 ms. We therefore further explore its underlying reasons in this subsection.

The cloud proxy on Gaikai is different from a local game console. It is a virtual machine (VM) running both computation-intensive tasks (for example, rendering the game) and bandwidth-intensive tasks (for example, streaming the game screen to the players) at the same time. Since these tasks cannot be decoupled into different VMs, it is necessary to see if they introduce extra overheads on the cloud proxy. Note that currently most cloud service providers usually provide the instances with custom optimization for DIA's specific needs [26] as well as the DIA's completely control over the instances [27], [28]. In other words, DIAs are usually assigned dedicated cloud instances. Thus, in this paper, we do not consider the interferences among different concurrent DIAs on the same instance.

Since the Limelight platform can hardly be tested by individual users, we run the standard CPU benchmark created by a tool called *sysbench* [29] on EC2 instances (their cloud platforms are quite similar in terms of Xen-based cloud virtualization). We use network tool *iPerf* [30] to adjust the traffic load on the instance and check the time cost of running the benchmark. To provide a fair comparison, we also do the experiment on local servers (non-virtualized servers) as a baseline. Figure 6 shows a comparison between an EC2 small instance and our local server. In this experiment, the EC2 small instance has 1.7 GB memory, 1 EC2 compute unit (Intel Xeon series, 1 virtual core with 1 EC2 compute unit), and 160 GB instance storage with 32-bit platform. Our local server also has similar hardware configuration that is comparable to the EC2 small instance. From this figure, we can see that the traffic load on the non-virtualized server only slightly increases the running time of the CPU benchmark, e.g., 250 Mbps traffic load will only increase the running time by 20%. However, for the virtualized EC2 small instance, this traffic load will double the running time of the CPU benchmark with very small standard deviation (the detailed

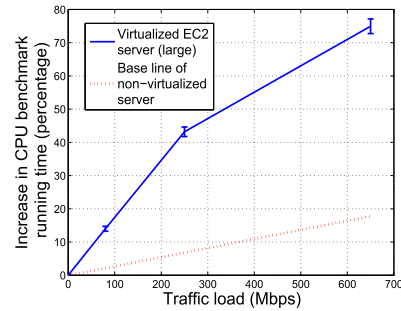


Fig. 7. Increase of CPU benchmark running time (on EC2 large instance).

TABLE II
CPU BENCHMARK RUNNING TIME UNDER DIFFERENT TRAFFIC LOADS ON VIRTUALIZED EC2 SMALL INSTANCE

Load	Run 1	Run 2	Run 3	Run 4	Avg
	Running time of CPU benchmark(ms)				
1 Mbps	37.17	36.71	36.75	36.64	36.82
5 Mbps	36.92	37.06	37.07	37.07	37.03
10 Mbps	37.41	37.51	37.36	37.42	37.42
22 Mbps	38.21	38.49	38.75	38.30	38.44
34 Mbps	39.03	39.13	39.06	39.13	39.09
40 Mbps	40.33	40.59	40.39	40.20	40.38
55 Mbps	41.28	41.22	41.06	41.67	41.31
75 Mbps	43.20	43.21	43.14	43.69	43.31
140 Mbps	49.00	48.84	50.06	48.13	49.01
190 Mbps	56.28	56.26	55.60	55.96	56.02
230 Mbps	71.50	70.67	69.06	73.01	71.06

data can be found in Table II). We have also tested this on EC2 large instances with multiple cores. The large instance has 7.5 GB memory, 4 EC2 computation units (Intel Xeon E5-2676 v3, 2 virtual cores and each with 2 EC2 computation units), 850 GB instance storage with 64-bit platform and very high I/O performance. Our local server, on the other hand, has weaker hardware configuration, particularly the CPU capacity. As shown in Figure 7, we can see that the traffic load on large instances still brings remarkable overheads to the system. Although the result is better than that of small instances, the traffic load will still remarkably slow down the running time of the CPU benchmark especially when comparing with the non-virtualized baseline.

It is easy to see that the CDIA cloud proxies are indeed in the same situation as these EC2 instances. The traffic load can significantly slow down the game running and unavoidably leads to a high processing latency. Yet, such a problem is rarely seen on the non-virtualized local game consoles or cloud proxies, or to a much lower degree.

V. LATENCY MINIMIZATION IN CDIA

Given the importance of interaction latency, there have been significant studies on latency minimization for conventional DIAs, mostly focusing on latency directly between client pairs [31] [32] [12]. Unfortunately, the existence of cloud

proxies prevents them from being used in CDIA, not to mention the task interference on the cloud proxies. In this section, we will revisit the latency modeling problem in these new contexts. We first consider a basic model to optimize the network latency. After that, this model will be further enhanced to capture the task interference on the cloud VMs.

A. Basic Model of CDIA Latency

To ensure responsive interactions, previous studies have suggested that reducing the average latency is not enough, because any fast users would suffer when they interact with long latency users [33] [34]. Our objective is thus to minimize the maximum latency between all client pairs that are bridged by cloud proxies. We focus on network latency here and will address processing latency in the next two sections.

$S = \{s_1, s_2, \dots, s_m\}$ denotes the set of m servers and $L = \{l_1, l_2, \dots, l_n\}$ denotes the set of n user clients. Let $C = \{c_1, c_2, \dots, c_o\}$ be the set of o cloud proxies. Each client in L will be assigned to a cloud proxy and a server in order to send operations and receive updates from other clients. An assignment A is a mapping from L to C and S , where for each client $l \in L$, $c_A(l) \in C$ denotes the assigned cloud proxy of client l and $s_A(l) \in S$ denotes the assigned server of client l .

For clients l_i and l_j , the communication should go through their assigned cloud proxies and servers in CDIA. Specifically, if l_i issues an operation to l_j , the following steps should be taken so that l_j can see the effect of this operation. First, l_i sends the operation to its assigned cloud proxy $c_A(l_i)$. $c_A(l_i)$ will then forward this operation to server $s_A(l_i)$ that is also assigned to l_i ; After that, if l_j is assigned to a different server $s_A(l_j)$, server $s_A(l_i)$ should forward the operation to server $s_A(l_j)$; Then $s_A(l_j)$ executes the operation and delivers the resultant state update to l_j 's cloud proxy $c_A(l_j)$; Finally, $c_A(l_j)$ will generate the game screen and stream the display to client l_j . Let $D(u, v)$ be the path latency between two nodes that are not directly connected and $d(u, v)$ be the link latency between two neighbor nodes. To be consistent with the existing DIA models [12], we assume that $D(u, v) = D(v, u)$ and $d(u, v) = d(v, u)$. The latency between client l_i and its server $s_A(l_i)$ can be calculated as:

$$D(l_i, s_A(l_i)) = d(l_i, c_A(l_i)) + d(c_A(l_i), s_A(l_i)) \quad (1)$$

We can therefore obtain the total interaction latency between l_i and l_j as follows:

$$D(l_i, l_j) = D(l_i, s_A(l_i)) + D(l_j, s_A(l_j)) + d(s_A(l_i), s_A(l_j)) \cdot I_{[s_A(l_i) \neq s_A(l_j)]} \quad (2)$$

where $d(s_A(l_i), s_A(l_j))$ denotes the latency between server $s_A(l_i)$ and $s_A(l_j)$, and $I_{[\cdot]}$ indicates whether l_i and l_j are assigned to different servers (1: yes; 0: no). Given the interaction latency between l_i and l_j , our objective is to find an assignment A to minimize $\mathbb{U}(A)$, the maximum interaction latency among all client pairs:

$$\text{minimize } \mathbb{U}(A) = \max_{l_i, l_j \in L} \{D(l_i, l_j)\} \quad (3)$$

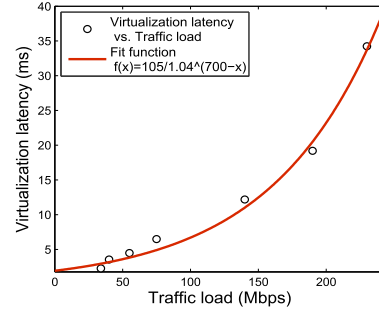


Fig. 8. Curve fitting for Equation (4) (on EC2 small instance).

This min-max latency can be optimally found when we convert it into a longest path problem in directed acyclic graph (DAG). The details are presented in *Appendix A* of the online supplementary material.

B. Enhanced Model to Capture Task Interference

In this part, we will further extend our model to consider the impact of traffic load on different cloud proxies. It is worth noting that CDIA offers elastic service capacity at cloud proxies. The cloud proxy capacity can be dynamically adjusted to meet user demands. Therefore, we denote P as the cloud proxy capacity set where $P = \{p_1, p_2, \dots, p_o\}$; $p_i \in P$ refers to the amount of resource that is assigned to cloud proxy c_i (bandwidth capacity in this case).⁷ Based on our measurement, we find that the NPV (Net Present Value) function [35] can be borrowed to capture the relationship between *virtualization latency* (processing latency caused by traffic load on VMs) and traffic load,⁸ we therefore compute the virtualization latency of cloud c_i as:

$$r(p_i) = \frac{a}{b^{p_i - q_A(c_i)}} \quad (4)$$

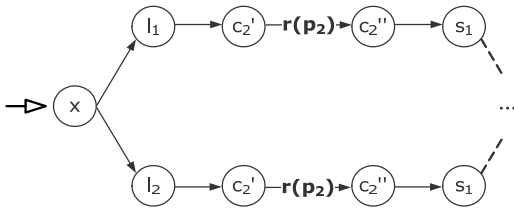
where a indicates the latency when the cloud proxy is fully loaded (without remaining bandwidth). Parameter b controls the skewness of the relationship between load and latency where $b \in (1, +\infty)$. Note that different VMs may have different a and b . For example, in Figure 8, a is around 105 and b is around 1.04.

Given a load assignment A and user l_i , we denote $p_A(l_i)$ as the resource assigned to cloud proxy $c_A(l_i)$. For a given set of servers, $S = \{s_1, s_2, \dots, s_m\}$ and clients $L = \{l_1, l_2, \dots, l_n\}$, the problem becomes how to use a set of cloud proxies $C = \{c_1, c_2, \dots, c_o\}$ to connect these clients and servers, with load assignment A and resource assignment P , to minimize the maximum interaction latency between all client pairs:

$$\begin{aligned} &\text{minimize } \mathbb{U}(A, P) \\ &= \text{Max}_{l_i, l_j \in L} \left\{ D(l_i, l_j) + r(p_A(l_i)) + r(p_A(l_j)) \right\} \quad (5) \end{aligned}$$

⁷We assume that one cloud proxy can serve multiple clients. Although using dedicated VMs reduces interference and improves the Quality of Experience (QoE), it dramatically increases the cost taking into account the large amount of users. Moreover, since we focus on the interference between computation-intensive tasks and bandwidth-intensive tasks, this interference still exists in the dedicated VM solution.

⁸This function has been widely used to quantify the relationship between cash and price/cost, which resembles our case when we try to purchase more cloud resources to reduce the virtualization cost on VMs.

Fig. 9. Transform G into G_A^* .

$$s.t. \forall i = 1, 2, \dots, o, \quad q_A(c_i) \leq p_i \quad (6)$$

$$\sum_{i=1}^o p_i * Cost(c_i) \leq K \quad (7)$$

where K refers to the total budget, which we assume can at least serve all the clients in the system.

It is easy to see that the virtualization latency makes the problem harder. If we assign client l_i to $c_A(l_i)$, it will not only assign traffic load to $c_A(l_i)$ but also affect the performance of other clients who have also been assigned to this cloud proxy. Assuming the cloud proxy capacity is given, this client assignment problem can therefore be transformed into a 0-1 Multiple Knapsack problem with a non-linear objective function, which is known to be NP-hard [36].

By exhaustively searching along all the possible combinations of A and P , the optimal solution can be achieved. However, the practical effectiveness of this search is limited considering the real-time user demands in CDIA systems. We thus propose a bi-level heuristic, which divides the optimization problem into two stages: load assignment and resource assignment. In the load assignment stage, we assume that all the cloud proxies are fully loaded (the virtualization latency is therefore equal to a in Equation (4)) and find the optimal load assignment A by Algorithm 5 (in Appendix A). After that, we construct a subgraph G_A based on the existing graph G and assignment A . As shown in Figure 9, we then split the node c_i into two virtual nodes (c_i' and c_i'') and use their link weight to refer to the virtualization cost on c_i . We denote G_A^* as the resulting graph and further apply a greedy algorithm to find the resource assignment P in G_A^* . As shown in Algorithm 1, this greedy algorithm iteratively assigns resource to the cloud proxies on the longest path. The algorithm stops when the remaining budget is not enough. In the next section, we will show that this bi-level heuristic achieves near-optimal performance in practical settings.

VI. PERFORMANCE EVALUATION

We now evaluate the performance of our solution via extensive trace-based simulations in MATLAB. The network latency (measured in Section IV) and the processing delay (measured in Section VI) will both serve as the inputs of our evaluation. We first examine the performance of our optimal client assignment when there is no task interference.⁹ After that, we investigate the performance of the interference-aware client assignment algorithm in the virtualized environment.

⁹This will be the case when the system is deployed on non-virtualized cloud platforms.

Algorithm 1 ResourceProvisioning()

```

1: Get  $G_A^*$  from  $A$ ;
2:  $R \leftarrow K - C$ ;
3: while (true) do
4:    $path^* = LongestPath(G_A^*)$ ;
5:   Get  $c_i, c_j$  from  $path^*$ ;
6:   if  $R \geq \max(Cost(c_i), Cost(c_j))$  then
7:     if  $\frac{r(p_i) - r(p_i + 1)}{Cost(c_i)} \geq \frac{r(p_j) - r(p_j + 1)}{Cost(c_j)}$  then
8:        $R \leftarrow R - Cost(c_i)$ ;
9:        $p_i \leftarrow p_i + 1$ ;
10:    else
11:       $R \leftarrow R - Cost(c_j)$ ;
12:       $p_j \leftarrow p_j + 1$ ;
13:    end if;
14:   else if  $R \geq \min(Cost(c_i), Cost(c_j))$  then
15:     if  $Cost(c_i) \leq Cost(c_j)$  then
16:        $w \leftarrow i$ ;
17:     else
18:        $w \leftarrow j$ ;
19:        $R \leftarrow R - Cost(c_w)$ ;
20:        $p(c_w) \leftarrow p(c_w) + 1$ ;
21:     end if;
22:   else
23:     break;
24:   end if;
25: end while;

```

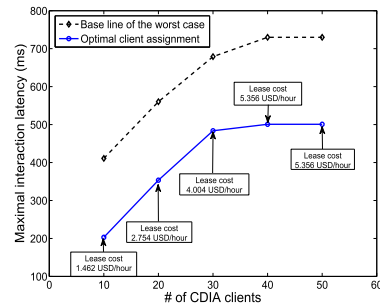


Fig. 10. Optimal client assignment only considering the networking latency.

We start with a CDIA system that consists of 20 clients, 5 cloud proxies and 5 servers. The renting cost of cloud proxies are referenced from the instance price list of Amazon's *On Demand instances* [16]. Figure 10 presents the performance of our optimal client assignment when there is no task interference (the processing latency is a default value of 80 ms at the cloud proxies). It is easy to see that the smart client assignment greatly reduces the interaction latency. Without optimization, the maximum client interaction latency can be as high as 700 ms. Our approach, on the other hand, can reduce the maximum latency to less than 500 ms. It is also worth noting that the renting price is linearly related to the client population. This indicates a good scalability of our approach.

It is not surprising to see that the optimal client assignment can achieve such a significant gain when there is no task interference. Figure 11 further explores the case when the optimal assignment can hardly be archived in the task

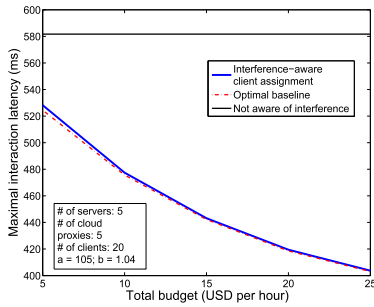


Fig. 11. Interference-aware client assignment.

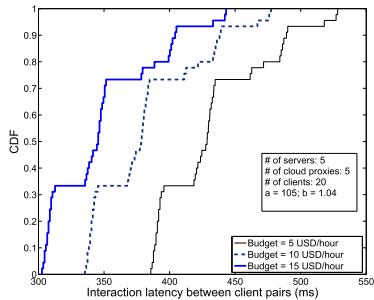


Fig. 12. Interaction latency across client pairs (with different budgets).

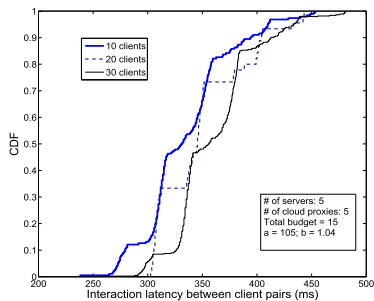
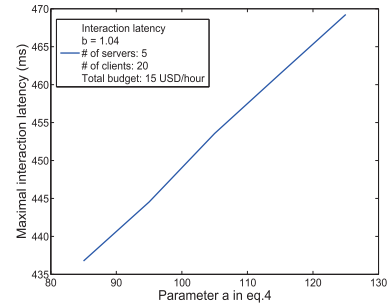
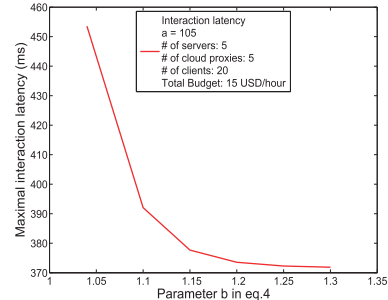


Fig. 13. Interaction latency across client pairs (different # of clients).

interference environment. We can see that task interference optimization is critical for CDIA. The maximum interaction latency can be larger than 580 ms if we focus on it only. Fortunately, our interference-aware algorithm can achieve a near-optimal (with the difference within 5 ms) latency that greatly reduces the interaction latency.¹⁰ It is worth noting that the interaction latency can be further reduced and become closer to the optimal results when we have more budget to purchase more capacities at the cloud proxies.

Figure 12 takes a closer look at the interaction latency between individual clients. We can see that all clients can benefit from the total budget increase. To be more specific, when the budget is equal to 5 USD/hour, less than 30% clients can have an interaction latency less than 400 ms. If we increase the budget to 15 USD/hour, more than 95% clients can interact with each other with a latency below 400 ms. The difference between the fastest and the slowest clients is also quite small, around 150 ms. Figure 13 further shows the cases

¹⁰The optimal base-line is obtained by brute-force searching.

Fig. 14. Adjusting parameter a (VM's maximum processing latency).Fig. 15. Adjusting parameter b (skewness of the relationship).

with different number of CDIA clients. We can see that for a given budget, our algorithm scales well with an increasing number of clients. Note that the total budget also bounds the total capacity of the cloud proxies. We thus cannot add more clients in Figure 13.

To understand the virtualization latency on different types of VMs, we investigate the case with different parameter inputs in Equation (4). Figure 14 presents the case when the maximum processing latency (parameter a for the cloud proxies) is changed from 85 ms to 125 ms. We can see that the interaction latency increases linearly with a . On the other hand, Figure 15 presents the case when virtualization latency and traffic load have more skewed relationships.¹¹ Based on these two figures, we can find that a good VM should have a small a and a large b . The maximum processing latency should be small when the VM is fully loaded (with a small a). In other words, adding idle resources on the VM ought to significantly reduce such a processing latency (with a large b).

Figure 16 further presents the CDF of the interaction latency across 200 random selected clients and all the 588 clients¹² in our measurement, respectively. It is easy to see that 80% clients can achieve the interaction latency within 300 ms. The interaction latencies between most (70%) client pairs are between 200 ms and 250 ms. It is also worth noting that the total budget in this case is relatively high up to 100 USD per hour. This is because we are using the pricing list of Amazon's *On Demand instances*. Choosing other types of platforms/instances, such as the *Reserved instance* may further reduce this cost.

¹¹Note that different a , b pairs in these two figures can present different cloud instances. For example, we use $a = 105$ and $b = 1.04$ to capture the features of EC2 small instances in our simulation.

¹²All the clients that we have used in our measurement in Section IV.

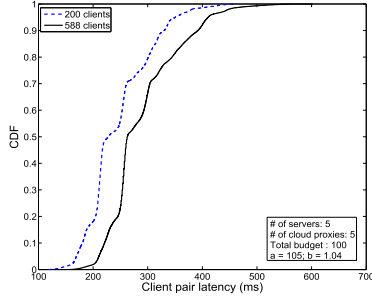


Fig. 16. Interaction latency across 200 and 588 clients.

VII. SYSTEM SCALING WITH USER DYNAMIC

The above model for the first time considers the task interference in classic CDIA systems. The related analysis can be applied to understand its budget-QoS relationship [37]. For example, how much money can archive what level of min-max interaction latency across how many users. In this section, we will further examine the system scaling issues with the arrival and the departure of clients. When a client arrives/departs the system, we need to minimize the maximum interaction latency with a bounded budget (in Equations (5), (6) and (7)). To better decide when we need to acquire/release extra cloud resource, we introduce a threshold ξ ($\xi \geq T_{max}$) to denote user's delay tolerance in different CDIA systems, where T_{max} is the maximum interaction latency across all clients. As we have discussed in Table I, the configuration of ξ is related to the type of different CDIA applications. This threshold adds the following new constraint to Equation (5):

$$\mathbb{U}(A, P) \leq \xi \quad (8)$$

Based on this new constraint, we then discuss our design to handle the client arrival/departure-aware step-by-step.

A. Handling Client Arrival and Departure

For any two clients $l_i, l_j \in L$, to emphasize the simultaneous consideration of communication and computation, we denote $T(l_i, l_j)$ as the interaction latency between client l_i and l_j , then we have $T_{max} = \max \{T(l_i, l_j)\}$. As to a new arriving client Δl , ΔT denotes the maximum interaction latency between Δl and the other clients, i.e., $\Delta T = \max \{T(\Delta l, l_i) | l_i \in L\}$. Since ΔT depends on the mapping from C to Δl , we define the min-max interaction latency across all proxies as $\Delta T_{min} = \min \{\Delta T(c_i) | c_i \in C\}$.

Based on Equations (5), (6), (7) and (8), Algorithm 2 gives the proxy assignment as well as the resource management approaches when a new client Δl joins the system. Note that in Step 2 of Algorithm 2, we calculate the expected processing time as

$$r^*(c_i) = \begin{cases} \frac{a}{b^{p_i - q(c_i)}}, & p_i = q(c_i) \\ \frac{a}{b^{p_i - q(c_i) - 1}}, & p_i > q(c_i) \end{cases} \quad (9)$$

The details of Steps 7 and 12 of Algorithm 2 are also given in Algorithm 3, in which CP denotes the cost performance of proxies. By expanding Equation (4) into a continuous function,

Algorithm 2 ClientArrival()

```

1: for all  $c_i$  that  $c_i \in C$  do
2:   Calculate the expected processing time  $r^*(c_i)$  when
    $c(\Delta l) \leftarrow c_i$ ;
3:    $\Delta T(c_i) = \max\{r^*(c_i) + D(\Delta l, l_j)\}, l_j \in L$ ;
4: end for;
5: Get  $\Delta T_{min} = \min \{\Delta T(c_i) | c_i \in C\}$ ;
6: if  $\Delta T_{min} > \xi$  then
7:   OnlineResourceProvisioning();
8: else
9:    $\Delta c = \{c_i | \Delta T(c_i) = \Delta T_{min}, c_i \in C\}$ ;
10:  Calculate  $T_{max}$  when  $c(\Delta l) = \Delta c$ ;
11:  if  $T_{max} > \xi$  then
12:    OnlineResourceProvisioning();
13:  else
14:     $c(\Delta l) \leftarrow \Delta c$ ;
15:    if  $p(c(\Delta l)) = q(\Delta l)$  then
16:       $p(c(\Delta l)) = p(c(\Delta l)) + 1$ ;
17:    end if;
18:  end if;
19: end if;
20:  $q(c(\Delta l)) = q(c(\Delta l)) + 1$ ;
21: return  $c(\Delta l)$ ;
```

we calculate CP as the differential coefficient of p_i , i.e.,

$$CP(c_i) = \text{diff} \left(\frac{r(p_i)}{p_i \cdot \text{Cost}(c_i)}, p_i \right) = \frac{-b^{q(c_i) - p_i} (1 + p_i \log b)}{p_i^2} \quad (10)$$

It is worth mentioning that when $\Delta T_{min} > \xi$, the interaction latency of the newly joined client becomes the bottleneck. In this case, we need to add more cloud resources on the proxies. To this end, we add Algorithm 3 to increase the resources on the proxies. This cyclic iterative algorithm is designed to clarify: 1) which proxy should be used to handle the new client, 2) whether we should add more resources to this proxy, and 3) how many resources should be added.

On the other hand, Algorithm 4 is designed to handle the client departure, where $path^*$ denotes the longest path before client Δl departs, T_{max} and T_{max}^* are the maximum interaction latency before and after client departure, respectively. We therefore summarize Theorem 2, of which the proof is detailed in Appendix D of the online supplementary material.

Note that according to Theorem 2, the client departure will not increase the maximum interaction latency. Hence, it is unnecessary to acquire proxy resources when considering the cost equality in Steps 6 and 10 of Algorithm 4.

B. Performance Validation

In this section, we aim to evaluate the performance of our resource provisioning design with user dynamic. As shown in Figure 17, we consider two typical patterns, *Stable* and *Flash-Crowd* [38] [39], in our evaluation. In particular, *Stable* means the number of clients is changing randomly between 30 to 50. *Flash-Crowd*, on the other hand, means the number of clients periodically increases from 5 to 50 and then decreases from

Algorithm 3 OnlineResourceProvisioning()

```

1:  $\Delta p \leftarrow 1$ ,  $target \leftarrow 0$ ,  $flag \leftarrow 1$ ;
2: Sort  $C$  by ascendant order of  $Cost(c_i)$ ;
3: while ( $flag$ ) do
4:   for all  $c_i$  that  $c_i \in C$  do
5:      $c(\Delta l) \leftarrow c_i$ ;
6:     Get the client  $l'$  that has the maximal interaction delay
       with  $\Delta l$ , then get  $c_j = c(l')$ ;
7:     if  $CP(c_i) > CP(c_j)$  then
8:        $p_j \leftarrow p_j + \Delta p$ ;
9:       Calculate  $T_{max}$ , i.e., the maximal interaction delay
       among all clients;
10:      if  $T_{max} \leq \xi$  then
11:         $target \leftarrow c_j$ ,  $flag \leftarrow 0$ ,  $c(\Delta l) \leftarrow c_j$ ;
12:      else
13:         $p_j \leftarrow p_j - \Delta p$ ;
14:      end if;
15:    end if;
16:    if ( $flag$ ) then
17:       $p_i \leftarrow p_i + \Delta p$ ;
18:      Calculate  $\Delta T(c_i)$ , i.e., the maximal interaction delay
       between  $\Delta l$  and the other clients;
19:      if  $\Delta T(c_i) \leq \xi$  then
20:         $target \leftarrow c_i$ ,  $flag \leftarrow 0$ ,  $c(\Delta l) \leftarrow c_i$ ;
21:      else
22:         $p_i \leftarrow p_i - \Delta p$ ;
23:      end if;
24:    end if;
25:  end for;
26:  if ( $flag$ ) then
27:     $\Delta p \leftarrow \Delta p + 1$ ;
28:  end if;
29: end while;
30: return  $c(\Delta l)$ ,  $\Delta p$ ,  $target$ ;

```

Algorithm 4 ClientDeparture()

```

1:  $q(c(\Delta l)) = q(c(\Delta l)) - 1$ ;
2:  $p(c(\Delta l)) = p(c(\Delta l)) - 1$ ;
3:  $path^* = LongestPath(G_A^*)$ ;
4: Get  $c_u, c_v$  from  $path^*$ ;
5: if  $CP(c_u) \geq CP(c_v)$  then
6:   if  $Cost(c_u) < Cost(c(\Delta l))$  then
7:      $p(c_u) = p(c_u) + 1$ ;
8:   end if;
9: else
10:  if  $Cost(c_v) < Cost(c(\Delta l))$  then
11:     $p(c_v) = p(c_v) + 1$ ;
12:  end if;
13: end if;

```

50 to 5 during a short period of time. These two patterns refer to typical “easy-to-handle” and “hard-to-handle” user demands, respectively.

Figure 18 shows the maximum interaction latency when user’s delay tolerance is set to 450 ms (a typical RPG game).

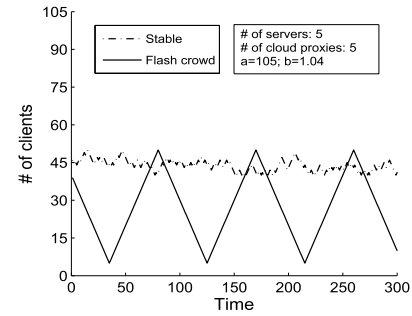
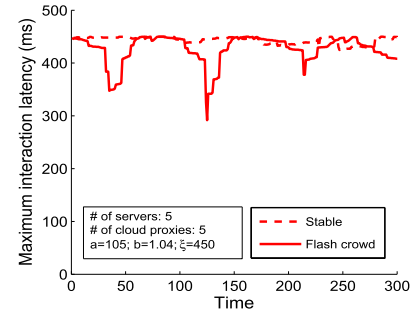
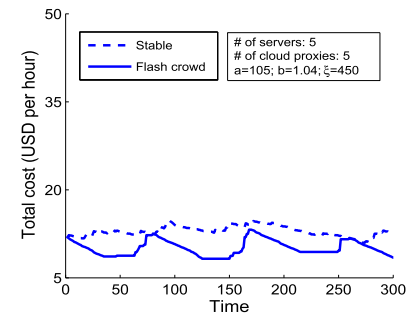


Fig. 17. Workload patterns.

Fig. 18. Maximum interaction latency over time ($\xi = 450$ ms).Fig. 19. Total cost over time ($\xi = 450$ ms).

As we can see in the figure, our proposed online algorithm can successfully bound the latency under both *Stable* and *Flash-Crowd* dynamics.¹³ Figures 19 and 20 further present the change of system cost. It is easy to see that the cost is less than 15 USD per hour (less than 3 cents/hour per user). Such a cost can be further reduced when we increase user’s delay tolerance. In particular, Figure 21 shows that the latency can be further reduced to less than 5 USD per hour (less than 1 cent/hour per user) when users’ delay tolerance is set to 650 ms. Although the interaction latency will be slightly increased, the maximum value is still bounded within 650 ms. As a summary, Figure 22 shows the relationship between cost and latency. With different configurations of ξ , our online protocol can always bound the maximum interaction latency with reasonable cost. It worth mentioning that since we dynamically manage resource when handling client arrival and

¹³The system’s maximum interaction latency is highly related to the number of clients in the system. Such jitter may also affect the QoE on different clients. The related QoE optimization is, however, out of the scope of this paper.

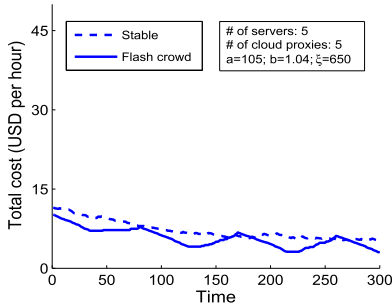
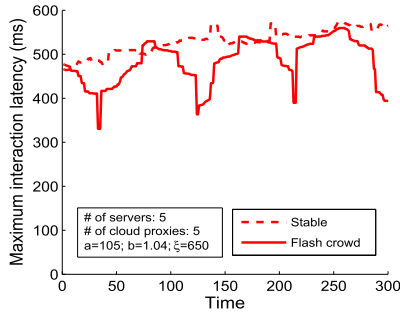
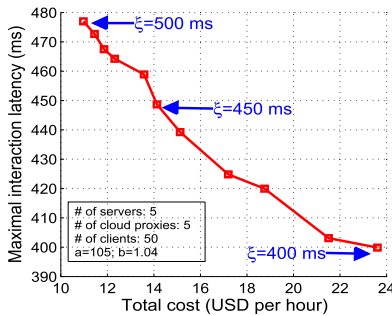
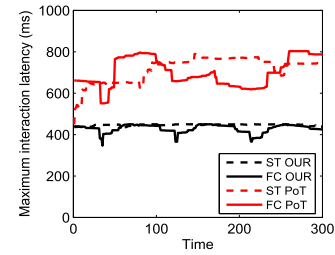
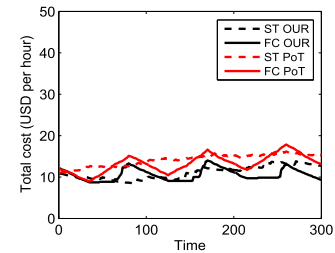
Fig. 20. Maximum interaction latency over time ($\xi = 650$ ms).Fig. 21. Total cost over time ($\xi = 650$ ms).

Fig. 22. Relationship between latency and cost.

departure, node state is updated in real time. Thus, when facing the fail-over (e.g., VM failure and migration), with the fast recovery capability of the cloud service provider, the proposed online protocol is capable to maintain the on-demand optimal allocations.

To illustrate the effectiveness, we compare the performance of our design and a Power-of-Two [40] strategy. Power-of-Two is a well-known solution to Supermarket Model, where in a supermarket a customer joins the server with shorter queue from two randomly sampled servers. Power-of-Two is an efficiency strategy, which can achieve a near-optimal waiting time for the customers. Because of its simplicity and effectiveness, Power-of-Two has been extensively researched [41]–[43] and widely used [44], [45]. In this simulation, Power-of-Two serves as a baseline. When a new client arrives, Power-of-Two randomly samples two cloud proxies, and then assigns the one with the smaller maximum interaction latency to the client. Meanwhile, Power-of-Two increases resource on the selected proxy. Under both stable (ST) and flash crowd (FC)

Fig. 23. Maximum interaction latency ($\xi = 450$).Fig. 24. Total cost ($\xi = 450$).

workload patterns when $\xi = 450$ ms, Figures 23 and 24 compare the maximum interaction latency and total cost of our solution (OUR) and Power-of-Two (PoT), respectively. As we can see, Power-of-Two induces larger maximum interaction latency as well as higher cost. Power-of-Two can hardly select the cloud proxy with minimum maximum interaction latency because of randomness. Besides, without considering the price/performance ratio when increasing resource, Power-of-Two tends to spend more money than our solution in operation.

VIII. RELATED WORK

The origin of Distributed Interactive Applications (DIAs) can be traced back to 1983 when a United States research program initiated the SIMNET project [46] to train soldiers in battlefield tactics. Since then, an increasing number of academic, military and commercial DIA systems have been developed and documented. Nowadays, despite the increase of processing powers at participating clients and the availability of greater communication bandwidth, minimizing the interaction latency remains one of the most fundamental challenges in the DIA framework. Many studies have shown that the latency is particularly problematic when the network delays are comparable to the interaction time or speed [47]. Such studies suggested that the interaction latency should be bounded for real-world DIAs [48]. For example, the typical latency values to maintain real-time interaction fall between 40 and 300 ms [49]. Gutwin [33] investigated the latency effects on two types of user interactions: movement prediction and moving a shared object. This study showed that both gaming performance and user strategy will be greatly affected by interaction latencies higher than the expected range.

To minimize the interaction latency in DIAs, Webb *et al.* [50] proposed a nearest server assignment to reduce the client-server latency. Ta and Zhou [51] proposed

a two-phase solution for large-scale DIAs. The study by Cronin *et al.* [31] further discussed the server placement problem to enhance user's interactivity. Vik *et al.* [32] explored the spanning tree problems in DIAs for latency reduction. Zhang and Tang [12] revisited the problem and proposed a distributed-modify-assignment approach to adapt to the dynamics of client participation and network conditions. A recent study from Lee *et al.* [52] presents *Outatime*, a speculative execution system for mobile cloud gaming that is able to mask up to 120 ms of network latency.

For cloud computing, there have been a series of works measuring the performance of public or private cloud services from diverse aspects, including computation, storage, and networking services [53]–[56]. There are also many studies addressing application designs that leverage cloud resources [57]–[63]. For example, Wu *et al.* [64] explored the use of cloud for Video-on-Demand applications [65], [66]. Huang *et al.* [67] provided an open-source cloud gaming system *GamingAnywhere*. Extensive experiments are performed to understand the video quality for both mobile and desktop clients. However, the deployment of enterprise cloud-based distributed interactive applications, which have emerged very recently, remains a mystery to the general public. Our study takes an initial step towards understanding this new generation of DIAs, in terms of its design, performance, and potential challenges.

IX. CONCLUSION AND FURTHER DISCUSSION

In this paper, we examined the framework design and latency optimization in cloud-based distributed interactive applications through real system measurement and analysis. Our study identified the unique features as well as the fundamental design challenges in CDIA. However, there are still many open issues that can be further explored in this new-born system.

First, to better mitigate the interference between the computation-intensive and bandwidth-intensive tasks, we are working on the analysis of TCP/UDP flows on different types of VMs. Our preliminary result shows that VM *hypervisors* (also known as *virtual machine managers* such as Xen, KVM and VMware) and VM total capacities play important roles for the interference. On the other hand, except for the Amazon EC2, other representative cloud service providers such as GCE [26], IBM Cloud [68], Rackspace [69], Microsoft Azure [70], and AliCloud [27] have also played an important role in promoting the rapid development of CDIA. Our future work also includes the comparison between different cloud service providers regarding network performance.

Second, we are currently investigating the efficiency by directly migrating some DIA protocols/optimizations into the CDIA framework. This analysis helps to better enjoy the benefits of cloud computing while minimize the corresponding overheads. These investigations not only improve the overall performance of the CDIA framework, it also helps to better understand the development of many other cloud-based systems with similar design frameworks.

Third, handling real-world scaling issues is also an important part of our future work. The solution proposed in this

paper requires global knowledge about the system state such as path latency and workload distribution. This may potentially affect the scalability of the solution. To mitigate this challenge, our on-going work is to build a hadoop2-like [71] [72] distributed information management models to collect global information, which enhances the efficiency and robustness of the system in a large-scale real-world environment.

ACKNOWLEDGMENT

Grateful acknowledgement is made to Dr. K. Yang and L. Lv who gave considerable help during experiments.

REFERENCES

- [1] *Essential Facts About the Computer and Video Game Industry*. Accessed: Oct. 1, 2013. [Online]. Available: http://www.theesa.com/facts/pdfs/ESA_EF_2012.pdf
- [2] *Entertainment Business*. Accessed: Oct. 1, 2013. [Online]. Available: <http://www.entertainmentbusiness.nl/sites/default/files/documents/2012/Videogames.pdf>
- [3] Gaikai. Accessed: May 1, 2017. [Online]. Available: <http://www.gaikai.com/>
- [4] Onlive. Accessed: May 1, 2017. [Online]. Available: <http://www.onlive.com/>
- [5] Ciinow. Accessed: May 1, 2017. [Online]. Available: <http://www.ciinow.com/>
- [6] *Gaikai in Play Station4*. Accessed: Oct. 1, 2013. [Online]. Available: <http://www.theverge.com/2013/2/20/4010420/sonys-playstation-4-will-use-gaikai-game-streaming-technology>
- [7] *Cloud-Based Gaming Planned for Microsoft XBOX 720*. Accessed: Oct. 1, 2013. [Online]. Available: <http://sonyps4playstation.com/cloud-based-gaming-planned-for-microsoft-xbox-720/>
- [8] *AMD Invests Into Ciinow*. Accessed: Oct. 1, 2013. [Online]. Available: <http://www.ubergizmo.com/2012/09/amd-invests-into-ciinow/>
- [9] Jupyter. Accessed: May 1, 2017. [Online]. Available: <http://jupyter.org/>
- [10] *Shiny: A Web Application Framework for R*. Accessed: May 1, 2017. [Online]. Available: <http://shiny.rstudio.com/>
- [11] F. Safaei, P. Boustead, C. D. Nguyen, J. Brun, and M. Dowlatshahi, "Latency-driven distribution: Infrastructure needs of participatory entertainment applications," *IEEE Commun. Mag.*, vol. 43, no. 5, pp. 106–112, May 2005.
- [12] L. Zhang and X. Tang, "Optimizing client assignment for enhancing interactivity in distributed interactive applications," *IEEE/ACM Trans. Netw.*, vol. 20, no. 6, pp. 1707–1720, Dec. 2012.
- [13] M. Claypool and K. Claypool, "Latency and player actions in online games," *Commun. ACM*, vol. 49, no. 11, pp. 40–45, 2006.
- [14] *Gaikai Powered Cloud-Based Gaming on Samsung Smart TVs*. Accessed: Jun. 5, 2012. [Online]. Available: <http://www.engadget.com/2012/06/05/gaikai-powered-cloud-gaming-coming-to-samsung-smart-tvs/>
- [15] S. A. Baset and H. G. Schulzrinne, "An analysis of the Skype peer-to-peer Internet telephony protocol," in *Proc. IEEE INFOCOM*, Apr. 2006, pp. 1–11.
- [16] Amazon EC2. Accessed: May 1, 2017. [Online]. Available: <http://aws.amazon.com/ec2/>
- [17] *Limelight Networks*. Accessed: May 1, 2017. [Online]. Available: <http://www.limelight.com/>
- [18] R. Dittner and D. Rule, *The Best Damn Server Virtualization Book Period*. Rockland, MA, USA: Syngress, 2007.
- [19] R. Jain and S. Routhier, "Packet trains—Measurements and a new model for computer network traffic," *IEEE J. Sel. Areas Commun.*, vol. 4, no. 6, pp. 986–995, Sep. 1986.
- [20] R. Kawahara, E. K. Lua, M. Uchida, S. Kamei, and H. Yoshino, "On the quality of triangle inequality violation aware routing overlay architecture," in *Proc. IEEE INFOCOM*, Apr. 2009, pp. 2761–2765.
- [21] *Networking and Traffic Control On Linux*. Accessed: May 1, 2017. [Online]. Available: <http://tcng.sourceforge.net/>
- [22] R. Shea, J. Liu, E. C.-H. Ngai, and Y. Cui, "Cloud gaming: Architecture and performance," *IEEE Netw.*, vol. 27, no. 4, pp. 16–21, Jul./Aug. 2013.

- [23] K.-T. Chen, Y.-C. Chang, H.-J. Hsu, D.-Y. Chen, C.-Y. Huang, and C.-H. Hsu, "On the quality of service of cloud gaming systems," *IEEE Trans. Multimedia*, vol. 16, no. 2, pp. 480–495, Feb. 2014.
- [24] M. Claypool, D. Finkel, A. Grant, and M. Solano, "On the performance of OnLive thin client games," *Multimedia Syst.*, vol. 20, no. 5, pp. 471–484, 2014.
- [25] M. Jarschel, D. Schlosser, S. Scheuring, and T. Hofbeld, "An evaluation of QoE in cloud gaming based on subjective tests," in *Proc. Innov. Mobile Internet Serv. Ubiquitous Comput.*, 2011, pp. 330–335.
- [26] GCE. Accessed: May 1, 2017. [Online]. Available: <https://cloud.google.com/compute/>
- [27] AliCloud. Accessed: May 1, 2017. [Online]. Available: <https://www.aliyun.com/>
- [28] Amazon Web Service. Accessed: May 1, 2017. [Online]. Available: <http://aws.amazon.com/>
- [29] Sysbench. Accessed: May 1, 2017. [Online]. Available: <https://github.com/akopytov/sysbench>
- [30] iPerf. Accessed: May 1, 2017. [Online]. Available: <https://iperf.fr/>
- [31] E. Cronin, S. Jamin, C. Jin, A. R. Kurc, D. Raz, and Y. Shavitt, "Constrained mirror placement on the Internet," *IEEE J. Sel. Areas Commun.*, vol. 20, no. 7, pp. 1369–1382, Sep. 2002.
- [32] K.-H. Vik, P. Halvorsen, and C. Griwodz, "Multicast tree diameter for dynamic distributed interactive applications," in *Proc. IEEE INFOCOM*, Apr. 2008, pp. 1597–1605.
- [33] C. Gutwin, "The effects of network delays on group work in real-time groupware," in *Proc. ECSCW*, 2001, pp. 299–318.
- [34] S. K. Barker and P. Shenoy, "Empirical evaluation of latency-sensitive application performance in the cloud," in *Proc. ACM MMSys*, 2010, pp. 35–46.
- [35] S. A. Ross, "Uses, abuses, and alternatives to the net-present-value rule," *JSTOR Financial Manage.*, vol. 24, no. 3, pp. 96–102, 1995.
- [36] C. Cotta and J. M. Troya, "A hybrid genetic algorithm for the 0-1 multiple knapsack problem," in *Artificial Neural Nets and Genetic Algorithms*. Vienna, Austria: Springer, 1998, pp. 250–254.
- [37] H. Wang, R. Shea, X. Ma, F. Wang, and J. Liu, "On design and performance of cloud-based distributed interactive applications," in *Proc. IEEE ICNP*, Oct. 2014, pp. 37–46.
- [38] M. Mao and M. Humphrey, "Auto-scaling to minimize cost and meet application deadlines in cloud workflows," in *Proc. ACM High Perform. Comput., Netw., Storage Anal.*, 2011, pp. 1–12.
- [39] K.-T. Chen, P. Huang, and C.-L. Lei, "Game traffic analysis: An MMORPG perspective," *Comput. Netw.*, vol. 50, no. 16, pp. 3002–3023, 2006.
- [40] M. Mitzenmacher, "The power of two choices in randomized load balancing," *IEEE Trans. Parallel Distrib. Syst.*, vol. 12, no. 10, pp. 1094–1104, Oct. 2001.
- [41] A. W. Richa, M. Mitzenmacher, and R. Sitaraman, "The power of two random choices: A survey of techniques and results," in *Combinatorial Optimization*, vol. 9. Berlin, Germany: Springer, 2001, pp. 255–304.
- [42] A. Mukhopadhyay and R. R. Mazumdar, "Analysis of randomized join-the-shortest-queue (JSQ) schemes in large heterogeneous processor-sharing systems," *IEEE Trans. Control Netw. Syst.*, vol. 3, no. 2, pp. 116–126, Jun. 2016.
- [43] A. Mukhopadhyay, A. Karthik, and R. R. Mazumdar, "Randomized assignment of jobs to servers in heterogeneous clusters of shared servers for low delay," *Stochastic Syst.*, vol. 6, no. 1, pp. 90–131, 2016.
- [44] X. Ren, G. Anantharayanan, A. Wierman, and M. Yu, "Hopper: Decentralized speculation-aware cluster scheduling at scale," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 45, no. 4, pp. 379–392, 2015.
- [45] S. Legtchenko *et al.*, "XFabric: A reconfigurable in-rack network for rack-scale computers," in *Proc. USENIX NSDI*, 2016, pp. 15–29.
- [46] SIMNET. Accessed: May 1, 2017. [Online]. Available: <http://en.wikipedia.org/wiki/SIMNET/>
- [47] P. M. Sharkey, M. D. Ryan, and D. J. Roberts, "A local perception filter for distributed virtual environments," in *Proc. IEEE Virtual Reality Annu. Int. Symp.*, Mar. 1998, pp. 242–249.
- [48] M. Mauve, J. Vogel, V. Hilt, and W. Effelsberg, "Local-lag and timewarp: Providing consistency for replicated continuous applications," *IEEE Trans. Multimedia*, vol. 6, no. 1, pp. 47–57, Feb. 2004.
- [49] C. Diot and L. Gautier, "A distributed architecture for multiplayer interactive applications on the Internet," *IEEE Netw.*, vol. 13, no. 4, pp. 6–15, Jul. 1999.
- [50] S. D. Webb, S. Soh, and W. Lau, "Enhanced mirrored servers for network games," in *Proc. ACM SIGCOMM NetGames*, 2007, pp. 117–122.
- [51] D. N. B. Ta and S. Zhou, "A two-phase approach to interactivity enhancement for large-scale distributed virtual environments," *Comput. Netw.*, vol. 51, no. 14, pp. 4131–4152, 2007.
- [52] K. Lee *et al.*, "Outatime: Using speculation to enable low-latency continuous interaction for mobile cloud gaming," in *Proc. ACM MobiSys*, 2015, pp. 151–165.
- [53] S. Garfinkel, "An evaluation of Amazon's grid computing services: EC2, S3, and SQS," Harvard Univ., Cambridge, MA, USA, Tech. Rep. TR-08-07, 2008.
- [54] Y. Cui, Z. Lai, X. Wang, N. Dai, and C. Miao, "QuickSync: Improving synchronization efficiency for mobile cloud storage services," *IEEE Trans. Mobile Comput.*, to be published, doi: 10.1109/TMC.2017.2693370.
- [55] Y. Cui *et al.*, "Software defined cooperative offloading for mobile cloudlets," *IEEE/ACM Trans. Netw.*, vol. 25, no. 3, pp. 1746–1760, Jun. 2017.
- [56] K. He *et al.*, "DeyPoS: Deduplicatable dynamic proof of storage for multi-user environments," *IEEE Trans. Comput.*, vol. 65, no. 12, pp. 3631–3645, Dec. 2016.
- [57] Y. Seung, T. Lam, L. E. Li, and T. Woo, "CloudFlex: Seamless scaling of enterprise applications into the cloud," in *Proc. IEEE INFOCOM*, Apr. 2011, pp. 211–215.
- [58] J. Wu, C. Yuen, N.-M. Cheung, J. Chen, and C. W. Chen, "Enabling adaptive high-frame-rate video streaming in mobile cloud gaming applications," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 25, no. 12, pp. 1988–2001, Dec. 2015.
- [59] S. Chen, Y. Shi, B. Hu, and M. Ai, "Mobility-driven networks (MDN): from evolution to visions of mobility management," *IEEE Netw.*, vol. 28, no. 4, pp. 66–73, Jul./Aug. 2014.
- [60] T. Li *et al.*, "On efficient offloading control in cloud radio access network with mobile edge computing," in *Proc. IEEE ICDCS*, Jul. 2017, pp. 2258–2263.
- [61] T. Li *et al.*, "Towards minimal tardiness of data-intensive applications in heterogeneous networks," in *Proc. IEEE ICCCN*, Aug. 2016, pp. 1–9.
- [62] K. Wang, K. Yang, H.-H. Chen, and L. Zhang, "Computation diversity in emerging networking paradigms," *IEEE Wireless Commun.*, vol. 24, no. 1, pp. 88–94, Feb. 2017.
- [63] K. Wang, K. Yang, and C. Magurawalage, "Joint energy minimization and resource allocation in C-RAN with mobile cloud," *IEEE Trans. Cloud Comput.*, to be published, doi: 10.1109/TCC.2016.2522439.
- [64] Y. Wu, C. Wu, B. Li, X. Qiu, and F. C. M. Lau, "CloudMedia: When cloud on demand meets video on demand," in *Proc. IEEE ICDCS*, Jun. 2011, pp. 268–277.
- [65] K. Xu *et al.*, "Modeling, analysis, and implementation of universal acceleration platform across online video sharing sites," *IEEE Trans. Serv. Comput.*, to be published, doi: 10.1109/TSC.2016.2565470.
- [66] K. He *et al.*, "Secure independent-update concise-expression access control for video on demand in cloud," *Inf. Sci.*, vol. 387, pp. 75–89, May 2017.
- [67] C.-Y. Huang, C.-H. Hsu, and K.-T. Chen, "GamingAnywhere: an open-source cloud gaming platform," *ACM SIG Multimedia Records*, vol. 7, no. 1, pp. 3–5, 2015.
- [68] IBM Cloud. Accessed: May 1, 2017. [Online]. Available: <http://www.ibm.com/cloud-computing/>
- [69] Rackspace. Accessed: May 1, 2017. [Online]. Available: <https://www.rackspace.com/>
- [70] Microsoft Azure. Accessed: May 1, 2017. [Online]. Available: <https://azure.microsoft.com/>
- [71] G. Turkington and G. Modena, *Learning Hadoop 2*. Birmingham, U.K.: Packt, 2015.
- [72] A. Sharma and S. Vyas, "Hadoop2 yarn," *Int. J. Comput. Sci.*, vol. 3, no. 9, pp. 18–26, 2015.
- [73] I. Katriel, L. Michel, and P. Van Hentenryck, "Maintaining longest paths incrementally," *Constraints*, vol. 10, no. 2, pp. 159–183, 2005.



Haiyang Wang received the Ph.D. degree from Simon Fraser University, Burnaby, BC, Canada, in 2013. He is currently an Assistant Professor with the Department of Computer Science, University of Minnesota at Duluth. His research interests include cloud computing, peer-to-peer networking, social networking, big data, and multimedia communications.



Tong Li received the B.S. degree from the School of Computer Science, Wuhan University, Wuhan, China, in 2012, and the Ph.D. degree from the Department of Computer Science and Technology, Tsinghua University, Beijing, China, in 2017. He is currently a Research Staff with Huawei Technologies Co., Ltd. His research interests include network protocols, cloud/edge computing, network virtualization, and resource management.



Ryan Shea received the Ph.D. degree in computer science from Simon Fraser University, Burnaby, BC, Canada, in 2016. He received the certificate in university teaching and learning from Simon Fraser University, where he is currently a University Research Associate in big data systems. His research interests include computer and network virtualization, and performance issues in cloud computing.



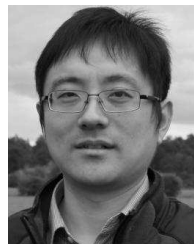
Xiaoqiang Ma received the B.S. degree from the Huazhong University of Science and Technology, Wuhan, China, in 2010, and the M.S. and Ph.D. degrees from Simon Fraser University, Burnaby, BC, Canada, in 2012 and 2015, respectively. He is currently an Assistant Professor with the School of Electronic Information and Communications, Huazhong University of Science and Technology. His current research interests include wireless networking, multimedia, cloud, and big data.



Feng Wang (S'07–M'13) received the B.S. and M.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, University of Mississippi, University, MS, USA. He was a recipient of the IEEE ICME Quality Reviewer Award in 2011. He is a Technical Committee Member of *Computer Communications*. He served as the Program Vice Chair for the International Conference on IOV 2014 and the TPC Co-Chair for the IEEE CloudCom 2017. He also serves as a TPC member in various international conferences, such as the IEEE INFOCOM, ICPP, IEEE/ACM IWQoS, ACM Multimedia, IEEE ICC, IEEE GLOBECOM, and IEEE ICME.



Jiangchuan Liu (S'01–M'03–SM'08–F'17) received the B.S. degree (*cum laude*) in computer science from Tsinghua University, Beijing, China, in 1999, and the Ph.D. degree in computer science from The Hong Kong University of Science and Technology in 2003. He was an Assistant Professor with The Chinese University of Hong Kong and a Research Fellow with Microsoft Research Asia. He is currently an EMC Endowed Visiting Chair Professor with Tsinghua University and an Adjunct Professor with the Tsinghua–UC Berkeley Shenzhen Institute. He is currently a University Professor with the School of Computing Science, Simon Fraser University, Burnaby, BC, Canada. His research interests include multimedia systems and networks, cloud computing, social networking, online gaming, big data computing, RFID, and Internet of Things. He is an NSERC E.W.R. Steacie Memorial Fellow. He was a co-recipient of the Inaugural Test of Time Paper Award of the IEEE INFOCOM in 2015, the ACM SIGMM TOMCCAP Nicolas D. Georganas Best Paper Award in 2013, and the ACM Multimedia Best Paper Award in 2012. He is a Steering Committee Member of the IEEE TRANSACTIONS ON MOBILE COMPUTING and the Steering Committee Chair of the IEEE/ACM IWQoS from 2015 to 2017. He has served on the Editorial Board of the IEEE/ACM TRANSACTIONS ON NETWORKING, the IEEE TRANSACTIONS ON BIG DATA, the IEEE TRANSACTIONS ON MULTIMEDIA, the IEEE COMMUNICATIONS SURVEYS AND TUTORIALS, and the IEEE INTERNET OF THINGS JOURNAL.



Ke Xu (M'02–SM'09) received the Ph.D. degree from the Department of Computer Science and Technology, Tsinghua University. He is currently a Full Professor with Tsinghua University. He holds three U.S. patents. His research interests include next-generation Internet, P2P systems, Internet of Things, network virtualization, and network economics. He is a member of the ACM. He received the Best Paper Award of GLOBECOM 2015. He serves as an Associate Editor of the IEEE INTERNET OF THINGS JOURNAL and has guest edited several special issues in IEEE and Springer journals.