

Immersion on the Edge: A Cooperative Framework for Mobile Immersive Computing

Zeqi Lai

Tsinghua University
laizq13@mails.tsinghua.edu.cn

Ziyi Wang

Tsinghua University
wangziyi0821@gmail.com

Yong Cui*

Tsinghua University
cuiyong@tsinghua.edu.cn

Xiaoyu Hu

Tsinghua University
chaose@gmail.com

CCS CONCEPTS

- **Networks** → *Cloud computing*;

KEYWORDS

Immersive Computing; Mobile Devices; Mobile Edge

ACM Reference Format:

Zeqi Lai, Yong Cui, Ziyi Wang, and Xiaoyu Hu. 2018. Immersion on the Edge: A Cooperative Framework for Mobile Immersive Computing. In *SIGCOMM Posters and Demos '18: ACM SIGCOMM 2018 Conference Posters and Demos, August 20–25, 2018, Budapest, Hungary*. ACM, New York, NY, USA, 3 pages. <https://doi.org/10.1145/3234200.3234201>

1 BACKGROUND AND MOTIVATION

1.1 Mobile IC apps are delay limited

Emerging *Immersive Computing (IC)* applications, such as virtual reality (VR) and augmented reality (AR), are changing the way human beings interact with mobile smart devices. It is well-known that *object recognition and rendering* are the key performance bottleneck in mobile IC systems [4, 6]. To speed up computation on mobile devices, the typical approach used in current IC applications is offloading computation-intensive tasks to the cloud [5, 6], or leveraging local system optimizations to accelerate IC tasks [2, 3]. However, as user's QoE requirements increase over time,

*Corresponding author. This project is supported by National Key R&D Program of China under Grant 2017YFB1010002, National 863 project (no. 2015AA015701).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGCOMM Posters and Demos '18, August 20–25, 2018, Budapest, Hungary

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-5915-3/18/08...\$15.00

<https://doi.org/10.1145/3234200.3234201>

future mobile IC applications demand higher visual quality (e.g., 4K or 8K resolution) and higher recognition accuracy/efficiency. Exploring innovations working with existing offloading approaches and local optimizations to enhance the performance of future mobile IC applications is still an important but challenging problem facing the IC industry.

1.2 Computation redundancy in IC apps

To gain insight on how to further reduce the user-perceived latency in modern IC applications, we analyzed more than 30 popular mobile VR/AR applications collected from Google Play and AppStore to understand the user interactions and computation workload. We derived three IC-specific insights, indicating that *IC tasks across different applications or users are often executed in similar or even redundant way*.

First, *object recognition* is a typical computation-intensive approach used in mobile AR or other vision-based assistant applications. We observe that it usually processes similar inputs in different applications/users. For example, two safe-driving applications are likely to recognize the same stop sign from the different angle at the same crossroads.

Second, we observe that under certain circumstance, interactive VR/AR applications require rendering the same 3D model on various devices. If multiple users play in the same environment, the content in the view of different users is likely to be similar. For example, two Pokemon Go players require rendering the same 3D avatar when they are interacting through Pokemon application in the same place.

In addition, current cloud-based VR applications leverage panoramic frames to create immersive experience [1, 4]. The server sends a panoramic frame to the client, and then the client crops the panorama to generate the final frame for display. Multiple users playing the same VR applications or watching the same VR video might use the same panorama.

In summary, we argue that computation-intensive tasks of mobile IC applications can be similar or redundant, especially when applications/users are in the close location. The above insights suggest the opportunity to improve QoE of

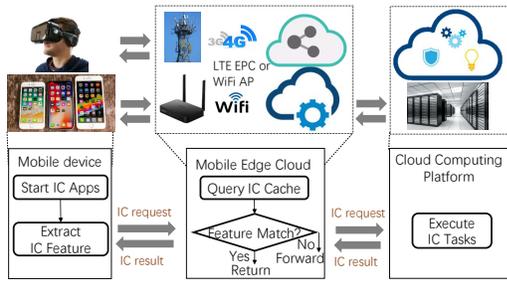


Figure 1: CoIC architecture.

immersive computing by cooperatively sharing and utilizing intermediate IC results among different applications/users.

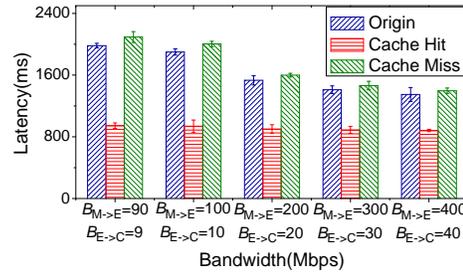
2 SYSTEM OVERVIEW

We present CoIC, a cooperative framework for mobile immersive computing applications. To speed up computation-intensive IC tasks, CoIC leverages the insight that similar or redundant IC tasks among different applications/users can be cached and shared to improve the user-perceived quality of experience (QoE), especially the end-to-end latency. Figure 1 shows the high-level architecture of CoIC. Initially, the client pre-processes the request to generate and send a *feature descriptor* of user’s input to the edge. On the edge, CoIC attempts to make a lookup with the feature descriptor (as the key) by matching the key to any results cached on the edge. If there is a hit, the cached result is returned to the client immediately. Otherwise, the edge forwards the request to the cloud and inserts the result to the edge cache.

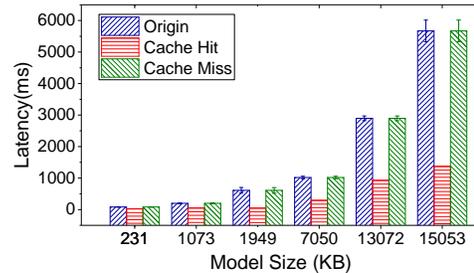
CoIC extracts dedicated property from each representative IC task as the feature descriptor. Specifically, for an object recognition task using DNN model, CoIC uses the feature vector generated from the input image as the feature descriptor. If the distance between the new feature descriptor and another one in the cache is under a certain threshold, CoIC determines that the computation result is already in the cache. For 3D object rendering and VR video streaming tasks, CoIC uses the hash value of the required 3D model or panoramic frames as the feature descriptor. Note that we did not cache object tracking results for AR applications because tracking is less computation-intensive as compared to recognition. Thus tracking is doable to be efficiently and accurately executed on mobile devices.

3 PRELIMINARY RESULTS

To evaluate the QoE improvement by CoIC, we implement an AR application upon CoIC, which renders high-quality 3D annotations to label objects recognized in the camera view. The client part of CoIC is implemented on a Pixel smartphone running Android O, and the edge/cloud modules of CoIC are implemented on two Linux machines respectively. The client connects to the edge via 802.11ac WiFi which supports up to 400Mbps available throughput in our experiment. We use `tc` to tune the network condition to simulate



(a) Recognition latency reduction under different network conditions. $B_{M \rightarrow E}$ and $B_{E \rightarrow C}$ refer to the available bandwidth between mobile client and edge, edge and cloud, respectively.



(b) Load latency reduction in rendering tasks.

Figure 2: Latency reduction of computation-intensive tasks (object recognition and 3D model loading).

real wireless/mobile network. In our experiment, the client sends recognition and rendering requests to the edge. The edge returns results immediately if the result is found in the cache. Otherwise the edge forwards the request to the server. We use an origin version which offloads complete IC tasks to the cloud without cache as the baseline. Figure 2a plots the reduction of recognition latency by CoIC. Current CoIC implementation recognizes objects via a DNN model. By identifying and caching similar computation results of the DNN model, CoIC can reduce up to 52.28% recognition latency under different network conditions. Figure 2b plots the latency reduction in rendering tasks. To execute a rendering task, the renderer has to load the 3D model into memory first and draw objects on the display. By caching the loaded data in rendering tasks on the edge, CoIC reduces the load latency by up to 75.86% for 3D models differed in size.

4 ONGOING AND FUTURE WORK

Our end goal is to improve the QoE of mobile IC applications by cooperatively utilizing the similar/redundant IC workload among different applications/users. Since the current CoIC can only identify coarse-grained IC tasks with simple cache management policy, we are exploring the improvement that can efficiently and accurately identify reusable IC workload in fine-grained (e.g., the result of a specific DNN layer). In addition, we will also study on the security/privacy protection issues in the cooperative system.

REFERENCES

- [1] Kevin Boos, David Chu, and Eduardo Cuervo. 2017. FlashBack: Immersive Virtual Reality on Mobile Devices via Rendering Memoization. In *MobiSys*. ACM, 23–27.
- [2] Peizhen Guo and Wenjun Hu. 2018. Potluck: Cross-Application Approximate Deduplication for Computation-Intensive Mobile Applications. In *ASPLOS*. ACM, 271–284.
- [3] Loc N Huynh, Youngki Lee, and Rajesh Krishna Balan. 2017. Deepmon: Mobile gpu-based deep learning framework for continuous vision applications. In *MobiSys*. ACM, 82–95.
- [4] Zeqi Lai, Y Charlie Hu, Yong Cui, Linhui Sun, and Ningwei Dai. 2017. Furion: Engineering High-Quality Immersive Virtual Reality on Today's Mobile Devices. In *MobiCom*. ACM, 409–421.
- [5] Xukan Ran, Haoliang Chen, Xiaodan Zhu, Zhenming Liu, and Jiasi Chen. 2018. DeepDecision: A Mobile Deep Learning Framework for Edge Video Analytics. In *INFOCOM*. IEEE.
- [6] Wenxiao Zhang, Bo Han, Pan Hui, Vijay Gopalakrishnan, Eric Zavesky, and Feng Qian. 2018. CARS: Collaborative Augmented Reality for Socialization. In *HotMobile*. ACM, 25–30.