# Toward Scalable and Controllable AR Experimentation

Ashkan Ganj
Worcester Polytechnic Institute
aganj@wpi.edu

Yiqin Zhao
Worcester Polytechnic Institute
yzhao11@wpi.edu

Federico Galbiati
Worcester Polytechnic Institute
fgalbiati@wpi.edu

Tian Guo
Worcester Polytechnic Institute
tian@wpi.edu

## ABSTRACT

To understand how well a proposed augmented reality (AR) solution works, existing papers often conducted tailored and isolated evaluations for specific AR tasks, e.g., depth or lighting estimation, and compared them to easy-to-setup baselines, either using datasets or resorting to time-consuming data capturing. Conceptually simple, it can be extremely difficult to evaluate an AR system fairly and in scale to understand its real-world performance. The difficulties arise for three key reasons: lack of control of the physical environment, the time-consuming data capturing, and the difficulties to reproduce baseline results.

This paper presents our design of an AR experimentation platform, ExpAR, aiming to provide scalable and controllable AR experimentation. ExpAR is envisioned to operate as a standalone deployment or a federated platform; in the latter case, AR researchers can contribute physical resources, including scene setup and capturing devices, and allow others to time share these resources. Our design centers around the generic sensing-understanding-rendering pipeline and is driven by the evaluation limitations observed in recent AR systems papers. We demonstrate the feasibility of this vision with a preliminary prototype and our preliminary evaluations suggest the importance of further investigating different device capabilities to stream in 30 FPS.

The ExpAR project site can be found at https://cake.wpi.edu/expar.

## CCS CONCEPTS

• **Human-centered computing** → **Ubiquitous and mobile computing design and evaluation methods**; • **Computer systems organization** → *Distributed architectures.*

## KEYWORDS

augmented reality, experimentation platform

## 1 INTRODUCTION

Augmented reality (AR) has emerged as a promising way for users to interact with physical worlds through virtual overlay. For example, in an AR-powered shopping app[1], a user can leverage a handheld device and its camera(s) to overlay products of interest, e.g., virtual glasses, on a desirable physical position, e.g., on the user's face [29]. As AR enters the general consumer market, hundreds of millions of users can benefit from this rich media experience with applications ranging from tourism to advertisement [6, 9].

Over the past decade, we have witnessed a blossom of works that provide high-quality AR performance [1, 10, 16, 20, 25, 28], including our work [30–32]. However, we have found that it is challenging to fairly and scalably evaluate the developed algorithms and systems to understand their real-world performance. We attribute the evaluation challenge to three key aspects: lack of control of the physical environment, time-consuming data capturing, and difficulty reproducing baseline results.

This paper introduces ExpAR, a platform providing scalable and controllable AR experimentation, by centering the key design insight of *generalizable AR pipelines*. Existing practices in capturing experimental data often involve a user interacting with the AR device [12, 26]. However, replicating the user mobility patterns, e.g., walking trajectory and device pose, and the environment information, e.g., scene lighting changes, is difficult. ExpAR controls the physical environment via programmable and remotely controllable mobility sensors, e.g., smart light bulbs and robotic cars [15].

Moreover, ExpAR provides mechanisms to capture high-quality data in a scalable and time-efficient manner. Access

---

[1]https://www.warbyparker.com/app

**Figure 1:** *An overview of* ExpAR *deployment.*

to high-quality data is crucial for properly evaluating an AR pipeline, from data captured by various sensors (most notably the camera sensors) to render virtual overlay. For example, prior work demonstrated that a blurred image could impact the accuracy of downstream vision tasks such as image classification and lighting estimation [16, 32]. ExpAR can capture high-quality data by supporting various capable hardware devices and carefully controlling their motions. Additionally, collecting a large amount of data in a scalable and time-efficient way is desirable, though such a desire is not unique to the AR community. For example, the robotic community has recently explored ways to enable multiple remote users to interact with robotic arms [22] and collect such data. ExpAR supports similar remote manipulation and allows programmable and parallel access to the physical capturing devices.

Lastly, to facilitate reproducibility, ExpAR breaks the AR pipeline into the sensing, understanding, and rendering steps. For each step, ExpAR will provide baseline methods that AR researchers can turn on and off to compose the desired AR pipeline dynamically. In other words, ExpAR will provide the ability to evaluate an AR solution *holistically* in the context of other AR tasks. To improve the baseline diversity, ExpAR will allow users to upload their own, similar to how Hugging Face hosts DL models [7].

Figure 1 depicts the high-level overview of ExpAR, a fully controllable and programmable AR evaluation platform. ExpAR is envisioned to operate as a standalone deployment (which we will describe an initial prototype in §4) or a federated platform that consists of network-connected deployments at different physical locations, similar to platforms such as PlanetLab and CloudLab [5, 14]. We envision ExpAR to consist of geographically-dispersed sites. Each site is a physical deployment that includes physical scene setups, AR, and capturing devices, that connect to ExpAR's backend for data storage and processing. Both AR researchers and users can interact with the physical setups remotely to carry out

key tasks, including scalable data capturing, experiment design, online surveys, and participant observation. We make the following key contributions.

- We pinpoint the limitations of existing AR evaluation methodology via characterizing recent papers and reflecting our evaluation practices.
- We describe the design of a fully controllable and programmable AR platform, centering the key insight of decomposable sensing-understanding-rendering AR pipelines. Our design serves as a conceptual framework for implementing a AR researcher-center evaluation platform.
- We present a preliminary prototype and evaluation that showcases the feasibility of programmable visual data capturing, streaming, and storage via a custom-built mobile capturing device and a cloud backend.

Our work shares similar spirits with three recent efforts, ILLIXR, XRBench, and CoMIC [11–13], in enabling better support for evaluating the emerging mixed reality applications. XRBench focuses on evaluating deep learning models for XR applications in representative execution patterns, which is similar to the design of ExpAR's *understanding* component. While ILLIXR and its multi-user counterpart CoMIC can capture data to evaluate a standalone system, it is not designed with a controllable physical environment, scalable experimentation, and cross-system evaluations in mind. In contrast, ExpAR is designed from outside to address the practical limitations exhibited in evaluating AR systems. In short, ExpAR compliments existing efforts and bridges the gap in reproducible AR research.

## 2 LIMITATIONS OF EXISTING AR EVALUATION METHODOLOGY

To understand the current practices and the limitations of AR evaluation, we surveyed 12 AR system papers focusing on their evaluation methodology. We categorize these papers based on AR tasks and characterize each paper's evaluation methodology along multiple important dimensions. The dimensions are selected based on our prior experiences in evaluating AR systems. Table 1 summarizes our findings where the *Task(s)* column refers to the number of tasks each paper evaluated. This analysis presents an in-depth understanding of the AR researcher's evaluation workflow. We make several key observations.

First, most works require capturing visual data during experimentation; some even used specialized hardware to gather ground truth [2, 10, 16, 24, 26]. Specialized hardware can be costly, e.g., Microsoft HoloLens 2 used by DeepMix costs $3.5K. Therefore, even if it is beneficial for evaluating AR systems on specialized hardware, not all papers can do so. However, if we can amortize the monetary cost by sharing the specialized devices among AR researchers, then it becomes more tractable to evaluate with specialized hardware.

**Table 1: *A survey of recent AR systems work and their evaluation methodology.*** For the last three columns, information inside the parenthesis represents the numerical scale. For example, Y(30) in the user study column means 30 participants.

| Category | Paper | Simulation | Task(s) | Visual data capturing (Y*- specialized) | Variation (**S**patial, **T**emporal) | Scene Diversity (**L**ow, **M**edium) | User Study | Eval. Scale (**S**mall, **M**edium) |
|---|---|---|---|---|---|---|---|---|
| **Lighting** | Gleam [20] | N | 1 | Y | S | L(2) | Y(30) | M(4) |
| | Xihe [31] | Y(Replay) | 1 | Y | T | L(1) | N | S(2) |
| | LitAR [32] | Y(Game Engine) | 1 | Y | S,T | L(3) | N | S(1) |
| **Depth** | InDepth [28] | Y(Dataset) | 1 | Y | S | M(3/20) | Y(27) | S(2) |
| | MobiDepth [26] | N | 2 | Y* | S | - | N | S(3) |
| **Tracking** | EdgeSLAM [2] | N | 1 | Y* | S,T | M(4) | N | S(2) |
| | AdaptSLAM [3] | Y | 1 | N(Dataset) | S,T | M(6) | N | S(1) |
| | FollowupAR [24] | N | 1 | Y* | S | M(2/variations) | N | S(1) |
| **Recognition/ Detection** | CollabAR [16] | Y(Dataset) | 1 | Y* | S | M | N | S(3) |
| | DeepMix [10] | Y(Dataset) | 1 | Y* | S | L(-/3) | Y(33) | S(1) |
| **Scheduling** | Heimdall [25] | N | K | Y | T | - | N | S(2) |
| **Multi-User** | SEAR [27] | Y | K | Y | T | - | N | S(1) |

This suggests the need for an experimentation platform to provide the visual data-capturing feature and time-sharing of expensive specialized hardware.

Second, many works evaluated and reported the performance with temporal and spatial variations. However, there is often a lack of *explicit control* of the physical environment. For example, in MobiDepth [26], the impact of spatial variance was measured by creating a dynamic scene that requires either moving the capturing device or the object of interest. However, the movement speed was only an approximated range, e.g., moving slowly vs. moving quickly. Although it was not explicitly mentioned, we suspect that the authors manually captured the required data. As such, it is hard to accurately quantify the impact of spatial locations on the AR systems (the proposed and the baselines). Temporal variations are slightly easier to control, e.g., in Xihe, we studied the impact of light intensity by fixing the rendering position and using a remotely controlled light source [31]. However, comparing how different systems work under the same temporal variations is non-trivial. We developed a session recorder to ensure consistent input to different systems [31].

Third, we saw that almost all papers have low scene diversity and small-scale evaluation scenarios. For example, Xihe was only evaluated in one physical room to understand its real-world performance [31]. Even for works with higher scene diversity, they were only evaluated in up to six scenes [3]. Additionally, most works evaluated use one to two mobile devices, with the upper end being four. Because of the heterogeneous mobile capabilities, it is hard to generalize and understand how well these works will perform in the wild. The requirements of having access to physical experiment spaces and high manual setup efforts seem to place a high toll on conducting diverse and large-scale evaluations.
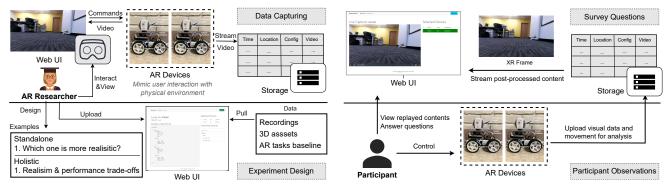
Fourth, the nature of AR research calls for visual perception studies. All papers presented metric-based qualitative

evaluations but only some conducted user studies to understand human perception performance. User studies are often considered to have a higher barrier for entry, e.g., requiring researchers to recruit and manage participants and design scalable study protocols.

Last, almost all papers focus on single-task evaluation. However, delivering the AR experience from sensory data to rendered results to end users involve a complex ecosystem and many moving pieces (see Figure 3). For example, an AR shopping app often requires hand tracking and object detection models to allow users to interact with the mixed environment [25]. Furthermore, in the context of supporting multi-user AR experiences [16, 27], it is unavoidable to consider inter-task dependencies. Evaluating one task in isolation is a good start, but we believe it would be better to evaluate the proposed solution in the context of the AR ecosystem. Such *holistic evaluations* can provide valuable insights into how well the proposed solution will work in a real-world deployment.

## 2.1 Reflections of Our Experimentation

We re-examine what we did when evaluating an AR system LitAR [32]. We used two groups of questions (Appendix B); the first group summarizes our evaluation process and rationales, while the second group asks for desirable features that can enable better evaluation experiences. We summarize the key reflection takeaways below.

**Controlled evaluation is time-consuming and difficult to set up.** Because AR evaluation often involves interaction with the physical environment, we need to control relevant physical factors during the experiments. For example, in LitAR, we need to control the *distance between observation and rendering positions* for each experiment run when evaluating its impact on lighting reconstruction quality. In the simplest case, this would involve manually setting up the capturing device and the props (i.e., a metal sphere ball) at specified locations and repeating the process for different

**(a)** AR researcher

**(b)** AR user study participant

**Figure 2: *An overview of* ExpAR *key workflows.*** AR researchers and user study participants can leverage ExpAR to perform key evaluation tasks, including data capturing, experiment design, participant observation, and online surveys.

distance variables. Each setup can take a few minutes, and therefore it can take many hours to complete a set of evaluations. However, this simple setup does not control other scene properties, e.g., lighting conditions and moving objects. Building a fully controlled physical scene is difficult (but not impossible). Instead, we resorted to a photorealistic indoor simulator that took about two months to set up.

**Evaluation diversity and scale are limited by monetary cost and time.** When it comes to input data diversity, we are often limited by whatever off-the-shelf sensory devices are available and the budget to acquire them. For example, it would be interesting to see how different LiDAR sensors impact LitAR's performance, but one such device (e.g., iPad Pro) costs more than $1K. Evaluating LitAR in more than three physical scene setups would be beneficial. Still, we were constrained by access to physical spaces and the ability to re-organize the scene (e.g., we don't have much control when using a public space). We also found ourselves developing *one-off* tools and workflows when capturing various sensory data, streaming these data to the edge, and managing these data for further analysis. Though this is a similar challenge to edge offloading, the added burdens of dealing with hardware sensors and different AR tasks make it non-trivial to scale up the data capturing and, thus, evaluations.

**Comparative studies are often guided by the easiness of reproducibility.** In LitAR, we only compared with two easy-to-setup baselines: a commercial solution ARKit and our prior work Xihe [31]. Would we benefit from comparing LitAR to other baselines? Probably. But such attempts were squashed by the hurdles to reproduce without source code and datasets or even the tremendous efforts required to set up the baselines. These hurdles also apply to user studies, which have other challenges, including participant recruitment [28] and multi-user coordination [22]. Additionally, we only evaluated how LitAR performs for the lighting estimation task; we did not evaluate LitAR in an AR application

to understand how it interacts with other AR tasks and its impact on the AR experiences. Even though this type of holistic evaluation is valuable for understanding in-the-wild performance, we see very few works that include holistic evaluations. We suspect that the lack of holistic evaluations is not caused by a lack of interest, but rather the lack of *plug-and-play* evaluation support. It is already difficult to set up baselines for a single task; we can't imagine the obstacles one has to overcome to configure an entire application scenario that requires the coordination of many tasks.

Besides these three key observations, we suspect the fundamental problem that limits the AR evaluation methodology is *treating evaluation as an afterthought*. In essence, we first design and build AR solutions and at a later stage, very reluctantly start the evaluation. The reluctance in part can be caused by the abovementioned challenges and obstacles in commencing any evaluations. ExpAR aims to simplify the evaluation process for AR solutions and promotes the *key principle of iteration, prototyping, and testing*.

## 3 EXPAR DESIGN

This section describes our overall vision of ExpAR, a fully controllable and programmable AR platform, and its high-level design goals. ExpAR is envisioned to operate as a standalone deployment (which we will describe an initial prototype in §4) or a federated platform that consists of network-connected deployments at different physical locations, similar to platforms such as PlanetLab and CloudLab [5, 14].

**Key design insight:** *let's generalize the AR pipeline!* Based on our experiences in designing AR systems and surveying other works (see Table 1), we present a generic *sensing-understanding-rendering* paradigm to which AR systems can generalize (Figure 3). Data capturing, i.e., *sensing*, plays a critical role in using deep learning models to *understand* the interaction between physical and virtual worlds. Those two components converge in the *rendering* in which virtual objects are composited and shown to the end users [4, 23]. By
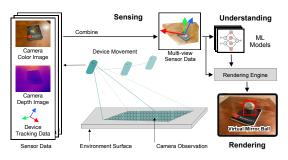
**Figure 3:** *The generic sensing-understanding-rendering pipeline for AR.* Sensing leverages the increasingly rich sensors, understanding relies on compute-intensive DL models, and rendering overlays augmented information using modern rendering engines.
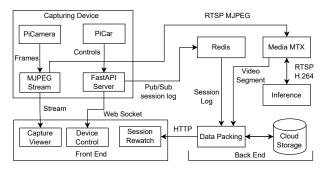
decomposing the AR pipeline into these three components, ExpAR can allow better sharing of each component among different active experiments and support holistic evaluations with minimal efforts from the AR researchers.

**Design goals.** We leverage our findings described in §2 and design ExpAR with the following three key goals. *(i) Controllable evaluation environment.* ExpAR should provide the ability to control each physical scene programmatically. This includes but is not limited to controlling physical environment conditions such as lighting and object placements and data capturing devices. *(ii) Scalable data capturing and parallel evaluations.* ExpAR should allow both AR researchers and users to access a wide variety of hardware devices, e.g., to capture data in parallel from devices residing in different physical locations. Devices can be time-shared, and different user study participants can use different devices to enable truly large-scale evaluations. *(iii) Reusable pipeline components and composable application scenarios.* ExpAR should provide built-in and default methods for different pipeline components. AR researchers can use existing components to define and configure their evaluation pipelines. To boost the component diversity, ExpAR will also allow community contribution, similar to existing AI/ML platforms like Hugging Face [7]. Moreover, ExpAR will provide default application templates that AR researchers can drag and drop their tasks into, as well as the ability to customize the templates for easy-to-setup holistic evaluations.

**Overview.** Figure 2 illustrates the key design components of a single deployment ExpAR and how two stakeholders, i.e., AR researchers and user study participants, interact with ExpAR. In a federated deployment, we will support a third stakeholder, the platform admin, who manages the operation of ExpAR. Additional example workflows are in Appendix A.

The basic setup of ExpAR consists of *data capturing devices*, a physical indoor scene where those capturing devices are initially parked and will explore, and a backend that persists data, including captured experimental data, static assets, and



**Figure 4:** *A prototype implementation.* The capturing device is based on Raspberry Pi 3, and the backend runs in Google Cloud as microservices inside docker containers.
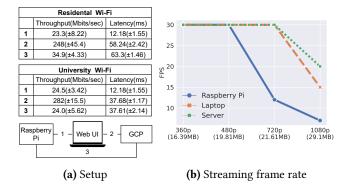
models/algorithms for different AR tasks. Many off-the-shelf hardware devices can act as the capturing devices, provided these devices are mobile, can be controlled remotely, and are equipped with necessary sensors, e.g., 360° RGB cameras and depth sensors. Example capturing devices include remote control cars and low-cost robots [15, 18, 19].

At a high level, AR researchers will programmatically control the capturing devices to traverse the physical scene and collect data helpful for understanding the environment, e.g., 360° videos. These data can be live viewed for monitoring and debugging purpose but will also be saved to facilitate future experiments. When an AR researcher needs to use ExpAR to capture the initial environment data, she can use any supported clients, e.g., VR headset or web UI, to select the desired physical scene. The scene configuration, which includes the number of capturing devices, capturing device capability, and the scene setup, like where the table is, is fixed for a given time. AR researchers select suitable scenes based on their experiment needs. Note that AR researchers do not need to have physical access to each scene and do not need to own any of the physical resources. AR researchers are the users of ExpAR instead of the owners. This design shares the same spirit as existing experimental platforms such as PlanetLab, CloudLab, and AWS's device farm [5, 8, 14].

Once researchers finish setting up the experiments, e.g., capturing required data and configuring the AR pipeline with desired template and components, user study participants can leverage ExpAR to conduct large-scale online surveys or remotely use the AR pipeline for participant observation.

## 4 PROTOTYPE IMPLEMENTATION

The current prototype of ExpAR is implemented in Python and Javascript in a containerized microservice architecture. Figure 4 presents an overview of ExpAR's key components and their interactions. The prototype consists of three logical modules: the front end, the capturing device, and the back end. The *front-end* module's primary objective is to provide ExpAR users, e.g., AR researchers and AR users, the ability to perform data capturing and monitor the capturing progress,

| Residental Wi-Fi | | |
|---|---|---|
| | Throughput(Mbits/sec) | Latency(ms) |
| 1 | 23.3(±8.22) | 12.18(±1.55) |
| 2 | 248(±45.4) | 58.24(±2.42) |
| 3 | 34.9(±4.33) | 63.3(±1.46) |

| University Wi-Fi | | |
|---|---|---|
| | Throughput(Mbits/sec) | Latency(ms) |
| 1 | 24.5(±3.42) | 12.18(±1.55) |
| 2 | 282(±15.5) | 37.68(±1.17) |
| 3 | 24.0(±5.62) | 37.61(±2.14) |

**(a)** Setup    **(b)** Streaming frame rate

**Figure 5: *Setup and streaming performance comparisons.*** Figure 5b shows ExpAR's FPS performance under the residential Wi-Fi[2].

as well as participate in the online survey via session rewatch. The *capturing device* module, encapsulating the hardware, is responsible for data collection and dispatching the collected data to the *back end*, which stores and processes the captured data. Currently, we implemented the front end in Vue.js, the data capturing based on a customized remote-controlled car PiCar-X [21] and a Raspberry Pi (RPi) 3B+, and the backend in containerized microservices running inside Google Cloud Platform (GCP)'s NVIDIA T4 GPU servers. See Appendix C for per-component implementations.

## 5 PRELIMINARY EVALUATION

We evaluate a prototype of ExpAR using a 360° camera, three representative computation devices, and a back end deployed to the Google Cloud Platform (GCP), to understand the streaming performance. The devices include: *(i)* a lower-end Raspberry Pi 3, equipped with a 1.2GHz 64-bit quad-core ARM Cortex-A57 CPU with 1GB RAM; *(ii)* a mid-end laptop, powered by an Intel Core i7-11370H processor with 16GB RAM; and *(iii)* a high-end server with a 4th Generation Intel Xeon processor and comes with 64GB RAM. Figure 5a depicts the setup and the pairwise network performance measured using `iPerf` under two Wi-Fi networks.

Figure 5b compares the frame per second (FPS) achieved under different streaming conditions. For each condition, we vary the frame resolution and the streaming device pairs. A streaming pair consists of one of the three devices (Raspberry Pi, laptop, and server) and an endpoint GCP server. Each device will stream the same video, recorded with the 360° camera, from a file. We make two key observations. First, as the resolution increases from 360p to 1080p, the FPS decreases for all devices. This suggests that the network starts to become the bottleneck (by comparing the network bandwidth to the streaming data size) and the device can't

---

[2]We observe similar FPS trends in university Wi-Fi.

**Table 2: *Task latency comparison over ten runs.***

| Task | Residential Wi-Fi | University Wi-Fi |
|---|---|---|
| Loading the system | 222.49 (± 8.55 ms) | 134.33 (± 3.93ms) |
| Setting up the device | 68.40 (± 8.42 ms) | 41.76 (± 2.53ms) |
| Client-server latency | 63.30 (±1.46 ms) | 38.06(± 2.24ms) |
| Executing control command | 6.11 (± 1.81 ms) | 1.81 (± 0.77ms) |

keep up with streaming larger 360° frames. Second, even the high-end server cannot achieve the desired 30 FPS when streaming at 1080p. We suspect that our CPU-based implementation is the culprit to such performance. We monitor the resource utilization and find that CPU utilizations increase drastically for larger resolutions (details in Appendix D). In short, our preliminary evaluations suggest the need to further investigate network optimizations and GPU-based implementations for supporting high-quality experimentation data streaming.

To better understand how ExpAR operates, we also evaluate the per-task latency under two Wi-Fi conditions. Table 2 summarizes the per-task latency: *(i) loading the system:* is the amount of time takes to get a response from our GCP Redis container. *(ii) setting up the device:* is the amount of time takes to write the device ID to the Redis database, *(iii) client-server latency:* describes the time to confirm that the server is active by pinging the server, *(iv) and executing commands:* is the amount of time takes for the PiCar-X to execute the user's commands. We find that it takes 176.09 ms and 290.89 ms to start up ExpAR using university and residential Wi-Fi respectively. We believe this overhead is reasonable because it is a one-time setup. The latency between client-server can be as high as 63.3 ms and can cause performance issues under certain interactive evaluation tasks. This suggests the need for a geo-distributed ExpAR deployment and considering the server locations based on ExpAR users' locations. University Wi-Fi in general leads to lower task latency and suggests the importance of properly configuring the network when using ExpAR for interactive experimentation.

## 6 CONCLUSION

This paper describes the design of ExpAR, an AR experimentation platform that allows setting up controllable evaluation environments easily, capturing data scalably, conducting evaluations in parallel, and reusing evaluation components. Our design is based on an in-depth analysis of the evaluation methodology from 12 recent AR system papers and our prior experiences and centers around the key insight of generalizable AR pipelines. ExpAR can allow AR researchers to share the physical devices and physical spaces, increasing the evaluation scale and diversity currently lacking. A prototype implementation and preliminary evaluation ExpAR were also presented, revealing interesting future directions such as improving the capturing device's onboard processing

power and optimizing the network performance between the front end and the capturing device. We will iterate the design and implementation of a local deployment while working on our AR projects. We hope to provide ExpAR as a service to the research community.

## ACKNOWLEDGMENTS

# REFERENCES

[1] Kittipat Apicharttrisorn, Jiasi Chen, Vyas Sekar, Anthony Rowe, and Srikanth V. Krishnamurthy. 2023. Breaking Edge Shackles: Infrastructure-Free Collaborative Mobile Augmented Reality *(SenSys '22)*. Association for Computing Machinery, New York, NY, USA, 1–15. https://doi.org/10.1145/3560905.3568546

[2] Ali J Ben Ali, Marziye Kouroshli, Sofiya Semenova, Zakieh Sadat Hashemifar, Steven Y Ko, and Karthik Dantu. 2022. Edge-SLAM: Edge-Assisted Visual Simultaneous Localization and Mapping. *ACM Trans. Embed. Comput. Syst.* 22, 1 (Oct. 2022), 1–31.

[3] Ying Chen, Hazer Inaltekin, and Maria Gorlatova. 2023. AdaptSLAM: Edge-Assisted Adaptive SLAM with Resource Constraints via Uncertainty Minimization. In *Proc. IEEE INFOCOM*.

[4] Ruofei Du, Eric Turner, Maksym Dzitsiuk, Luca Prasso, Ivo Duarte, Jason Dourgarian, Joao Afonso, Jose Pascoal, Josh Gladstone, Nuno Cruces, Shahram Izadi, Adarsh Kowdle, Konstantine Tsotsos, and David Kim. 2020. DepthLab: Real-time 3D Interaction with Depth Maps for Mobile Augmented Reality. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology (UIST '20)*.

[5] Dmitry Duplyakin, Robert Ricci, Aleksander Maricq, Gary Wong, Jonathon Duerig, Eric Eide, Leigh Stoller, Mike Hibler, David Johnson, Kirk Webb, Aditya Akella, Kuangching Wang, Glenn Ricart, Larry Landweber, Chip Elliott, Michael Zink, Emmanuel Cecchet, Snigdhaswin Kar, and Prabodh Mishra. 2019. The Design and Operation of CloudLab. In *Proceedings of the USENIX Annual Technical Conference (ATC)*. 1–14.

[6] The Green Planet AR Experience. 2022. https://www.factory42.uk/greenplanetexperience.

[7] Hugging Face. 2023. https://huggingface.co/.

[8] AWS Device Farm. 2023. https://aws.amazon.com/device-farm/.

[9] Immersive Stream for XR (Preview). 2023. https://xr.withgoogle.com/.

[10] Yongjie Guan, Xueyu Hou, Nan Wu, Bo Han, and Tao Han. 2022. DeepMix: mobility-aware, lightweight, and hybrid 3D object detection for headsets. In *Proceedings of the 20th Annual International Conference on Mobile Systems, Applications and Services (MobiSys '22)*. 28–41.

[11] Bo Han, Parth Pathak, Songqing Chen, and Lap-Fai Craig Yu. 2022. CoMIC: A Collaborative Mobile Immersive Computing Infrastructure for Conducting Multi-user XR Research. *IEEE Netw.* (2022), 1–9.

[12] Muhammad Huzaifa, Rishi Desai, Samuel Grayson, Xutao Jiang, Ying Jing, Jae Lee, Fang Lu, Yihan Pang, Joseph Ravichandran, Finn Sinclair, Boyuan Tian, Hengzhi Yuan, Jeffrey Zhang, and Sarita V Adve. 2021. ILLIXR: Enabling End-to-End Extended Reality Research. In *2021 IEEE International Symposium on Workload Characterization (IISWC)*. 24–38.

[13] Hyoukjun Kwon, Krishnakumar Nair, Jamin Seo, Jason Yik, Debabrata Mohapatra, Dongyuan Zhan, Jinook Song, Peter Capak, Peizhao Zhang, Peter Vajda, Colby Banbury, Mark Mazumder, Liangzhen Lai, Ashish Sirasao, Tushar Krishna, Harshit Khaitan, Vikas Chandra, and Vijay Janapa Reddi. 2023. XRBench: An extended reality (XR) machine learning benchmark suite for the metaverse. In *Proceedings of Machine Learning and Systems*, Vol. 5.

[14] Larry L. Peterson and Prof. David Culler. 2002. PlanetLab. https://planetlab.cs.princeton.edu/.

[15] Liangkai Liu, Ren Zhong, Aaron Willcock, Nathan Fisher, and Weisong Shi. 2023. An Open Approach to Energy-Efficient Autonomous Mobile Robots citation. In *2023 International Conference on Robotics and Automation (ICRA)*. IEEE Press, 7 pages.

[16] Z Liu, G Lan, J Stojkovic, Y Zhang, C Joe-Wong, and M Gorlatova. 2020. CollabAR: Edge-assisted Collaborative Image Recognition for Mobile Augmented Reality. In *2020 19th ACM/IEEE International Conference on Information Processing in Sensor Networks (IPSN)*. 301–312.

[17] Ajay Mandlekar, Yuke Zhu, Animesh Garg, Jonathan Booher, Max Spero, Albert Tung, Julian Gao, John Emmons, Anchit Gupta, Emre Orbay, Silvio Savarese, and Li Fei-Fei. 2018. ROBOTURK: A Crowd-sourcing Platform for Robotic Skill Learning through Imitation. In *Conference on Robot Learning*.

[18] Adithyavairavan Murali, Tao Chen, Kalyan Vasudev Alwala, Dhiraj Gandhi, Lerrel Pinto, Saurabh Gupta, and Abhinav Gupta. 2019. PyRobot: An Open-source Robotics Framework for Research and Benchmarking. *arXiv preprint arXiv:1906.08236* (2019).

[19] Oliver Kroemer. 2023. LoCoBot: An Open Source Low Cost Robot. http://www.locobot.org/.

[20] Siddhant Prakash, Alireza Bahremand, Linda D Nguyen, and Robert LiKamWa. 2019. GLEAM: An Illumination Estimation Framework for Real-time Photorealistic Augmented Reality on Mobile Devices. In *Proceedings of the 17th Annual International Conference on Mobile Systems, Applications, and Services* (Seoul, Republic of Korea) *(MobiSys '19)*. Association for Computing Machinery, New York, NY, USA.

[21] SunFounder. 2023. Raspberry Pi Ai Car Kit - PiCar-X. https://www.sunfounder.com/products/picar-x. Accessed: 2023-6-15.

[22] Albert Tung, Josiah Wong, Ajay Mandlekar, Roberto Martín-Martín, Yuke Zhu, Li Fei-Fei, and Silvio Savarese. 2021. Learning Multi-Arm Manipulation Through Collaborative Teleoperation. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*. 9212–9219.

[23] Jamie Watson, Mohamed Sayed, Zawar Qureshi, Gabriel J. Brostow, Sara Vicente, Oisin Mac Aodha, and Michael Firman. 2023. Virtual Occlusions Through Implicit Depth. In *CVPR*.

[24] Jingao Xu, Guoxuan Chi, Zheng Yang, Danyang Li, Qian Zhang, Qiang Ma, and Xin Miao. 2021. FollowUpAR: enabling follow-up effects in mobile AR applications. In *Proceedings of the 19th Annual International Conference on Mobile Systems, Applications, and Services (MobiSys '21)*. Association for Computing Machinery, 1–13.

[25] Juheon Yi and Youngki Lee. 2020. Heimdall: mobile GPU coordination platform for augmented reality applications. In *Proceedings of the 26th Annual International Conference on Mobile Computing and Networking* (London, United Kingdom) *(MobiCom '20, Article 35)*. Association for Computing Machinery, New York, NY, USA, 1–14.

[26] Jinrui Zhang, Huan Yang, Ju Ren, Deyu Zhang, Bangwen He, Ting Cao, Yuanchun Li, Yaoxue Zhang, and Yunxin Liu. 2022. MobiDepth: realtime depth estimation using on-device dual cameras. In *Proceedings of the 28th Annual International Conference on Mobile Computing And Networking (MobiCom '22)*.

[27] Wenxiao Zhang, Bo Han, and Pan Hui. 2022. SEAR: Scaling Experiences in Multi-user Augmented Reality. *IEEE Trans. Vis. Comput. Graph.* 28, 5 (May 2022), 1982–1992.

[28] Yunfan Zhang, Tim Scargill, Ashutosh Vaishnav, Gopika Premsankar, Mario Di Francesco, and Maria Gorlatova. 2022. InDepth: Real-time Depth Inpainting for Mobile Augmented Reality. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 1 (March 2022), 1–25.

[29] Yiqin Zhao, Sean Fanello, and Tian Guo. 2023. Multi-Camera Lighting Estimation for Photorealistic Front-Facing Mobile Augmented Reality. In *Proceedings of the 24th International Workshop on Mobile Computing Systems and Applications* (Newport Beach, California) *(HotMobile '23)*. Association for Computing Machinery, New York, NY, USA, 68–73. https://doi.org/10.1145/3572864.3580337

[30] Yiqin Zhao and Tian Guo. 2020. PointAR: Efficient Lighting Estimation for Mobile Augmented Reality. In *The European Conference on Computer Vision (ECCV '20)*. 678–693.

[31] Yiqin Zhao and Tian Guo. 2021. Xihe: a 3D vision-based lighting estimation framework for mobile augmented reality. In *Proceedings of the 19th Annual International Conference on Mobile Systems, Applications, and Services (MobiSys '21)*. 28–40.

[32] Yiqin Zhao, Chongyang Ma, Haibin Huang, and Tian Guo. 2022. LITAR: Visually Coherent Lighting for Mobile Augmented Reality. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 3 (Sept. 2022), 1–29.

# A   EXPAR EXAMPLE WORKFLOWS

## A.1   Design for AR Researchers

We envision that ExpAR can aid AR developers and researchers in evaluating AR systems better by providing features including *data capturing*, *sensing visualization*, *experimentation design*, and *reproducible pipeline evaluation*. This section describes the design of ExpAR by explaining how an AR researcher will use ExpAR to accomplish the key evaluation tasks of *data capturing* and *experimentation design*. Figure 2a depicts the workflows associated with these two evaluation tasks.

*A.1.1   Data capturing.* At the center of the *data capturing* task lies the capturing devices. We envision many capturing device types available in the ExpAR, which the AR researchers can choose from. For example, if the AR researcher is interested in capturing the ground truth for the lighting estimation task, she can select the capturing device that consists of 360° cameras. Note that the AR researcher does not need to be physically coupled with the capturing devices. Rather, the AR researcher will use any provided ExpAR client-side interfaces, e.g., a web UI, to control the movement of the capturing devices to mimic how an AR user will interact with the physical environment.

Moving the capturing device to different physical locations will allow the AR researcher to capture desired experimental data. These data, e.g., in the form of RGB video and movement commands, will be stored for later use. The AR researcher can start the data capturing *in parallel* with different devices as the physical resource permits. To provide full control of the physical scene data, we will only allow at most one AR researcher to use the physical deployment. This resembles the current experimental practice where researchers are in charge of setting up the scene and introducing known dynamics to the scene. In other words, we will impose the time sharing at the physical scene level, and the researcher can capture at most $\sum_i^m n_i$ stream of data in parallel where $m$ is the number of idle physical scenes, and $n_i$ is the number of capturing devices of $i^{th}$ scene. In addition to the parallel capturing, we will also allow researchers to monitor the capturing progress to spot any abnormalities, e.g., due to misconfigurations. In short, an AR researcher can use ExpAR to capture large amounts of experimental data from diverse scenes with heterogeneous devices to suit their evaluation needs. Even better, the data capturing can be done with minimal human effort and does not require the costly acquisition of specialized devices.

*A.1.2   Experiment design.* AR researchers can use ExpAR to setup the evaluation pipeline so that user study participants can use this pipeline to answer survey questions. Often the process involves AR researchers leveraging tools to create and generate visual assets for the survey questions and online survey platforms such as Quatrics to distribute the surveys. In other words, the survey question preparation and distribution are done in separate pipelines. This current practice works but can be time-consuming for AR researchers to design the survey questions. Instead, our goal in designing ExpAR is to streamline the experiment design process by allowing AR researchers to "plug-and-play" different pipeline components to assemble final visual products used in the survey questions. For example, AR researchers can take a raw video stream and pass it through an end-to-end pipeline to generate a post-processed video stream in which they can directly embed questions to suitable frame locations.

More concretely, with ExpAR, the researcher will come up with a list of questions, and use these questions to guide the configuration of the evaluation pipeline. For example, if the researcher is interested to evaluate the performance of her rendering-related AR tasks such as lighting and depth estimation, she can choose the *baseline methods* hosted by ExpAR to render virtual assets, e.g., racing cars, in the same physical scene capture. With the rendered results, she can then create survey questions that ask participants to compare the relative rendering performance between her proposed method and a baseline method. Because the researcher is evaluating an AR task in isolation, we refer to this type of evaluation as *standalone evaluation*.

Researchers can also perform *holistic evaluation* to study how well the proposed method works with the remaining sensing-perception-rendering pipeline. In holistic evaluations, ExpAR will ask researchers to select from predefined AR application scenarios such as museum tours or furniture shopping or configure their custom scenario. Each scenario specifies the AR tasks that need to be activated in the pipeline; for example, in the furniture shopping scenario [25], ExpAR will activate the baseline models for image segmentation, object detection, and hand tracking tasks. For custom scenarios, the researchers are responsible for specifying interested AR tasks and their execution dependency. Afterward, researchers can leverage ExpAR to compare the end-to-end performance and quality trade-offs between the proposed method and any baselines.

ExpAR will provide relevant data such as physical scene recordings (either directly recorded by this researcher or shared by others), 3D assets, and baselines for AR tasks. In addition, ExpAR will also allow researchers to upload any custom data and configure the data's visibility, private or public. Based on the evaluation questions, these data will be assembled *on-demand*, providing the flexibility to conduct plug-and-play AR evaluations in scale.

## A.2 Design for AR User Study Participants

This section describes the design of ExpAR by explaining how a user study participant will use ExpAR to complete two types of common evaluations: survey questions and participant observations.

*A.2.1 Participant observation.* It is valuable to understand how end users interact with a new AR system. However, it can be troublesome and time-consuming to invite user study participants to a physical experimentation site. Other extraordinary conditions such as the COVID-19 pandemic can also limit the practice of this experimentation form. For multi-user collaborative AR systems, it can be even more challenging to invite multiple AR users to be physically present in the experimentation site simultaneously. Furthermore, requiring physical presence also restricts the demographic diversity of the participants as people who are physically close by, e.g., university students, are more likely to participate in the study.

ExpAR aims to facilitate the participant observation experiments by relaxing the physical presence requirement. That is, AR users can remotely interact with the AR systems (and devices). This concept is similar to RoboTurk, a recent robotic framework that allows human users to demonstrate how to perform tasks [17]. Because of the relaxed physical presence requirement, ExpAR can then support more geographically diverse user study participants and make multi-user experiments easy to conduct. To allow AR researchers to observe the interactions, ExpAR will provide a web portal that in real-time displays AR users' interactions, an over-the-head view of the physical scene and the AR devices' movement, and the video streams from individual AR devices. Moreover, ExpAR will store those data in the backend to support post-analysis.

In short, ExpAR will need to support streaming device perception from a physical scene to where the AR user locates and then send the AR user's interactions back to the physical devices. This communication pattern is similar to cloud gaming, and therefore we suspect the challenges lie in designing network optimizations to provide low-latency interactions.

*A.2.2 Online survey.* Understanding how a human user perceives AR features is important, and such understanding is often achieved via user studies. A common way to conduct user studies [20, 28] is to invite participants to answer survey questions, e.g., how two competing techniques compare visually. One of the key features ExpAR can support is to allow user study participants to take these surveys via the web portal. As described in §A.1.2, these survey questions are designed ahead of time by AR researchers who will then invite participants by sharing the survey URLs. Depending on the survey questions, the participants might be watching a replayed video stream in which frames are overlayed with information such as environmental conditions and visual outputs from different AR models. Relevant survey questions will be displayed to the participants at pre-specified frames and responses will be collected, similar to how MOOC quizzes online learners' understanding of course topics.

## B REFLECTION QUESTIONS

### B.1 Group One

- *Describe the general process you took when designing and running the experiments in one of your recent AR works.*
- *How long did setting up the physical testbed take?*
- *How long did data capturing take?*
- *How long did obtaining results from baselines take?*
- *How many mobile devices did you evaluate your system on?*
- *What were some reasons that prevent you from evaluating your system on more mobile devices?*
- *What are the task(s) in evaluating the AR systems that you find yourself repeatedly doing all the time?*
- *Was this experience described above, the same as other research projects you have done? The same as other AR research projects?*

### B.2 Group Two

- *If you could have access to a magical experiment setup that would allow you to do anything you want to evaluate your system, what would this experiment setup look like?*
- *How would it work?*
- *What would you use this magical setup to do?*
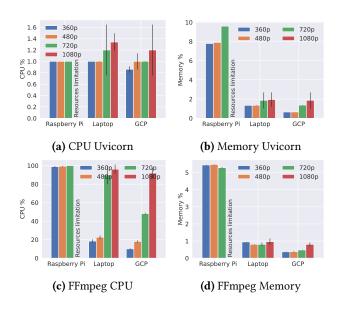
## C MODULE IMPLEMENTATION DETAILS

### C.1 Front end

We implemented a web UI to support the two stakeholders of ExpAR. Currently, it consists of three components: *(i)* The *Capture Viewer* component allows the AR researchers to watch a live stream from the capture device in the data collection process. *(ii)* The *Session Rewatch* component fetches the processed frames from the cloud storage and presents them to the user for user studies, allowing assessment and analysis of previously recorded AR sessions. *(iii)* The *Device Control* component processes controls from researchers, e.g., in the form of configuration files, and then controls various data capturing devices to initiate the data capture process. In the future, we also plan to support other UIs, such as VR headsets and controllers.

### C.2 Capturing Device

This module interfaces with the hardware sensors and the onboard computational resources. Currently, the hardware consists of a 360° camera, a PiCar-X, and an RPi 3B+. Our modular design allows for integrating additional capturing

**(a)** CPU Uvicorn



**(b)** Memory Uvicorn



**(c)** FFmpeg CPU



**(d)** FFmpeg Memory

**Figure 6:** *Resource utilization comparison.* We measure both CPU and memory utilization when streaming using residential Wi-Fi. We see that FFmpeg consumes significantly higher CPU as resolution increases.

devices, such as drones and mobility-enabled specialized devices, in the future.

We implemented two key modules: the *FastAPI server* for processing user control commands and *FFmpeg* for streaming video frames in MJPEG to the front/back ends. We have FFmpeg send the incremental frame counter to the FastAPI server to synchronize the time between controls and the video frames. Specifically, the streaming functionality was implemented by having a process to read PiCamera frames directly from the camera and then send them to the backend through FFmpeg via an RTSP connection executed in a sub-process. We implemented a continuous frame counter that increments with each new frame captured and shared the counter across all tasks.

We also used FastAPI to create a WebSocket connection between the front end and the capturing device. Once this connection is established, the live video feed will be activated and can be streamed to the back end. Moreover, the Web-Socket connection remains open and actively listens for user inputs. These inputs, formatted as device control instructions, are forwarded to the device interface for execution.

### C.3 Back end

Our back end consists of four main components for storing, processing, and streaming the captured data. The back end was implemented as Docker microservices, making the setup easily reproducible. The back end processes the inbound data

flow from the capturing device, including video, audio, and metadata.

The *MediaMTX* server encodes the raw MJPEG stream into H.264, a format that can be streamed to the front end. It also generates the video segments for the *data packing* component and interacts with the *inference* service to augment the video stream. The *Redis* component logs user actions, which will be supplied to the *data packing* component.

The *Data Packing* component post-processes and consolidates all the data recorded during an AR session into a single MP4 file, at the end of each session. It runs a Redis pub/sub subscribe listener and will start a background process when messages are published to the data packing channel. It retrieves events from the Redis stream log and encodes them as JSON text to generate a SubRip Text (SRT) file. We use SRT, rather than KLV, to include synchronized metadata in the streams because KLV does not have good open-source support. It then uses FFmpeg to combine video segments and the SRT file into a single MP4 file. Inference output data can also be encoded as subtitles.

The *Inference* component runs DL models on the video stream. Currently, we implemented a popular object detection model called YOLOv8 as a proof-of-concept. Upon initiation by MediaMTX, it receives the video stream and applies the YOLOv8 model to each video frame. It returns the video to MediaMTX, which subsequently streams it back to the client and stores it for future playback purposes. Our modularized design makes integrating other AR models into ExpAR as microservices easier.

## D RESOURCE UTILIZATION RESULTS

To better understand the performance bottlenecks of ExpAR, we measured the resource utilizations under different hardware setups. To facilitate these tests, we developed a mock car script that emulates car movements. This enables us to capture and stream content on the three computation devices while simultaneously measuring system resource usage.

As shown in Figure 6, higher video quality consumes more resources, especially by the FFmpeg process, which is responsible for stream and video format conversion. The Raspberry Pi struggled with high-quality video (1080p) processing, with FFmpeg maxing out CPU usage, and caused watchdog reset. Our results on the laptop and cloud server had a similar trend, with the CPU utilization increasing with the resolution. For higher resolutions, FFmpeg is again experiencing CPU bottlenecks. This indicates that the current implementation of ExpAR is not effectively utilizing system resources. Raspberry Pi's performance limitations and underutilization of resources in more powerful hardware highlight future improvement areas in exploring techniques to improve resource utilization and ensure smooth, high-quality AR experimentation experiences on a wide array of hardware.