

HoloLens 2 Research Mode as a Tool for Computer Vision Research

Dorin Ungureanu Federica Bogo Silvano Galliani Pooja Sama Xin Duan
 Casey Meekhof Jan Stühmer* Thomas J. Cashman Bugra Tekin
 Johannes L. Schönberger Pawel Olszta Marc Pollefeys
 Microsoft

Abstract

Mixed reality headsets, such as the Microsoft HoloLens 2, are powerful sensing devices with integrated compute capabilities, which makes it an ideal platform for computer vision research. In this technical report, we present HoloLens 2 Research Mode, an API and a set of tools enabling access to the raw sensor streams. We provide an overview of the API and explain how it can be used to build mixed reality applications based on processing sensor data. We also show how to combine the Research Mode sensor data with the built-in eye and hand tracking capabilities provided by HoloLens 2. By releasing the Research Mode API and a set of open-source tools, we aim to foster further research in the fields of computer vision as well as robotics and encourage contributions from the research community.

1. Introduction

Mixed reality technologies have tremendous potential to fundamentally transform the way we interact with our environment, with other people, and the physical world. Mixed reality headsets such as the Microsoft HoloLens 1 & 2 have already seen great adoption in a number of fields, particularly in first-line worker scenarios, ranging from assisted surgery to remote collaboration and from task guidance to overlaying digital twin on the real world. Despite existing adoption in these fields, the general space of mixed reality is still in its infancy. Oftentimes, the development of new mixed reality applications requires fundamental research and the novel combination of different sensors. The entry barrier to computer vision research in this area is significantly lowered by access to tools enabling us to effectively collect raw sensor data and develop new computer vision algorithms that run on-device.

The first-generation HoloLens Research Mode, released in 2018, enabled computer vision research on device by pro-

*Now at Samsung AI Centre, Cambridge (UK). Work performed while at Microsoft.



Figure 1. Microsoft HoloLens 2.

viding access to all raw image sensor streams – including depth and IR. Released together with a public repository collecting auxiliary tools and sample applications [4], Research Mode promoted the use of HoloLens as a powerful tool for doing research in computer vision and robotics [1].

HoloLens 2 (Fig. 1), announced in 2019, brings a number of improvements with respect to the first-generation device – like a dedicated DNN core, articulated hand tracking and eye gaze tracking [2]. However, being a novel platform built on new hardware, HoloLens 2 is not compatible with the previous version of Research Mode.

In this technical report, we introduce the second-generation HoloLens Research Mode, made available in 2020. Research Mode provides a set of C++ APIs and tools to access the HoloLens 2 sensor streams. We discuss the main novelties with respect to the previous version and present a set of applications built on top of it. For additional material and API documentation, we refer the reader to the GitHub repository [3].

Research Mode is designed for academic and industrial researchers exploring new ideas in the fields of computer vision and robotics. It is not intended for applications deployed to end-users. Additionally, Microsoft does not provide assurances that Research Mode will be supported in future hardware or OS updates.

The rest of this technical report is organized as follows.

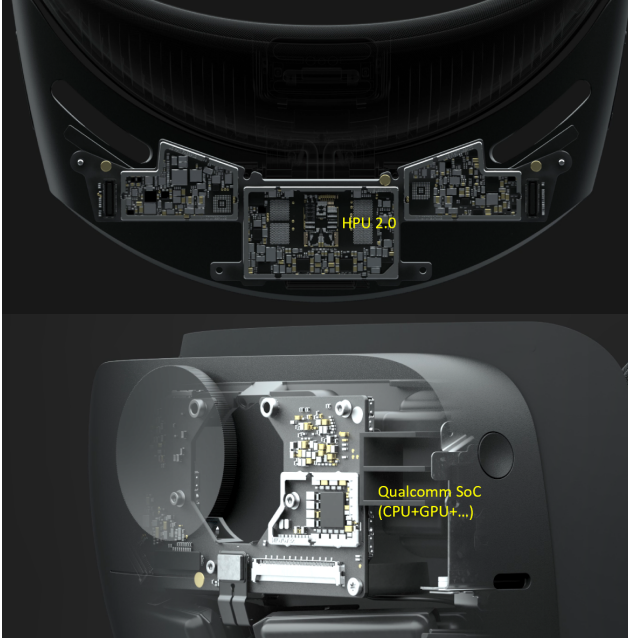


Figure 2. HPU and SoC: The HPU (top) is located on the frontal part of the device, near the sensors. The SoC (bottom), a Qualcomm Snapdragon 850, is located on the rear.

Section 2 introduces the HoloLens 2 device, detailing in particular its input streams. Section 3 provides an overview of the Research Mode API, while Section 4 showcases a few example applications. Finally, Section 5 sums up our contributions.

2. HoloLens 2

The HoloLens 2 mixed reality headset brings a set of improvements with respect to the first-generation device – including a larger field of view, a custom DNN core, fully articulated hand tracking and eye gaze tracking.

The device features a second-generation custom-built Holographic Processing Unit (HPU 2.0), which enables low-power, real-time computer vision. The HPU runs all the computer vision algorithms on device (head tracking, hand tracking, eye gaze tracking, spatial mapping etc.) and hosts the DNN core. It is located on the front part of the device, near the sensors (Fig. 2, top). The CPU on the SoC (a Qualcomm Snapdragon 850) remains fully available for applications. The SoC is located on the rear (Fig. 2, bottom).

The device is equipped with a depth and an RGB camera, four grayscale cameras, and an Inertial Measurement Unit (IMU), as shown in Fig. 3. Audio is captured with a microphone array (5 channels).

Research Mode for HoloLens 2 enables access to the following input streams:

Stream	Resolution	Format
VLC	640×480	8-bit
Long throw depth	320×288	16-bit
Long throw AB	320×288	16-bit
AHAT	512×512	16-bit
AHAT AB	512×512	16-bit

Table 1. Research Mode camera frames.

- Four visible-light tracking cameras (VLC): Grayscale cameras (30 fps) used by the system for real-time visual-inertial SLAM.
- A depth camera, which operates in two modes:
 - AHAT (Articulated HAnd Tracking), high-framerate (45 fps) near-depth sensing used for hand tracking. As hands are supported up to 1 meter from the device, the HoloLens 2 saves power by calculating only “aliased depth” from the phase-based time of flight camera. This means that the signal contains only the fractional part of the distance from the device when expressed in meters (see Fig. 4).
 - Long Throw, low-framerate (1-5 fps) far-depth sensing used to compute spatial mapping on device.
- Two depth modes of the IR stream (Active Brightness, AB in short), computed from the same modulated IR signal for depth computation. These images are illuminated by infrared and unaffected by ambient visible light (see Fig. 4, right).
- Inertial Measurement Unit (IMU):
 - Accelerometer, used by the system to determine the linear acceleration along the x , y and z axes as well as gravity.
 - Gyroscope, used by the system to determine rotations.
 - Magnetometer, used by the system for absolute orientation estimation.

For each stream, Research Mode provides interfaces to retrieve frames and associated information (*e.g.*, resolution and timestamps), and to map the sensors with respect to the device and to the world. Table 1 summarizes the main characteristics of each (camera) input stream. In the following section, we provide an overview of the API.

3. Research Mode API

In this section, we describe the main objects exposed by Research Mode and how they are used to access the sensor

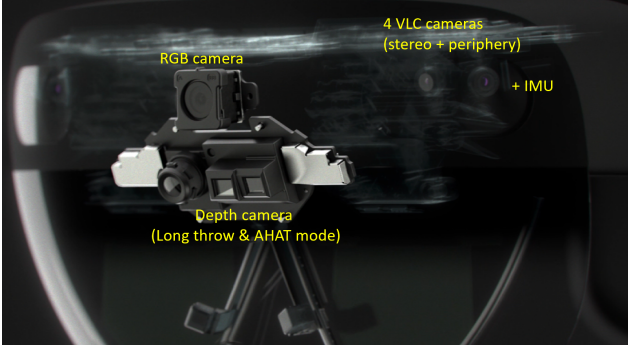


Figure 3. HoloLens 2 input sensors.

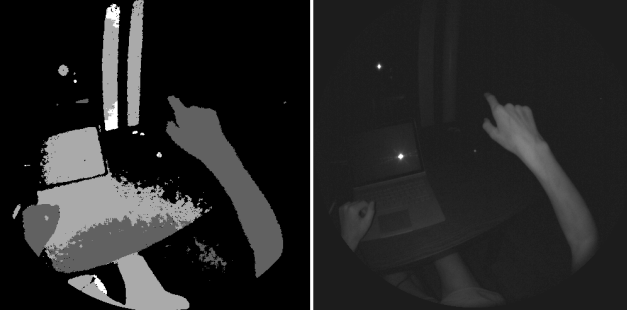


Figure 4. Depth in AHAT mode: Depth (left) and Active Brightness (right) images.

input streams. We refer the reader to [3] for a detailed API documentation.

The Research Mode main sensor loop (Sec. 3.1) starts by creating a `ResearchModeDevice` object, which is used to obtain a list of available `Sensors`. `Sensor` objects expose methods to retrieve and process frames (Sec. 3.2) and to locate the sensors with respect to the device and the world (Sec. 3.3).

3.1. Main Sensor Loop

The main sensor processing loop involves instantiating a `ResearchModeDevice`, getting sensor descriptors, opening sensor streams and fetching frames:

```

HRESULT hr = S_OK;
IResearchModeSensorDevice *pSensorDevice;
std::vector<ResearchModeSensorDescriptor>
    sensorDescriptors;
size_t sensorCount = 0;

hr = CreateResearchModeSensorDevice(&pSensorDevice);

pSensorDevice->DisableEyeSelection();

hr = pSensorDevice->GetSensorCount(&sensorCount);
sensorDescriptors.resize(sensorCount);
hr = pSensorDevice->GetSensorDescriptors(
    sensorDescriptors.data(),
    sensorDescriptors.size(),
    &sensorCount);

```

```

for (auto& sensorDescriptor : sensorDescriptors)
{
    // Sensor frame read thread
    IResearchModeSensor *pSensor = nullptr;
    size_t sampleBufferSize;
    IResearchModeSensorFrame* pSensorFrame = nullptr;

    hr = pSensorDevice->GetSensor(
        sensorDescriptor.sensorType,
        &pSensor);
    hr = pSensor->GetSampleBufferSize(&sampleBufferSize);
    hr = pSensor->OpenStream();

    for (UINT i = 0; i < 4; i++)
    {
        hr = pSensor->GetNextBuffer(&pSensorFrame);
        if (pSensor->GetSensorType() >= IMU_ACCEL)
        {
            // Process IMU frame
            SaveFrameImu(pSensor, pSensorFrame, i);
        }
        else
        {
            // Process camera frame
            SaveFrameCamera(pSensor, pSensorFrame, i);
        }
        if (pSensorFrame)
            pSensorFrame->Release();
    }
    hr = pSensor->CloseStream();
    if (pSensor)
        pSensor->Release();
}
pSensorDevice->EnableEyeSelection();
pSensorDevice->Release();
return hr;

```

Note that `OpenStream` and `GetNextBuffer` need to be called from the same thread. Since `GetNextBuffer` calls are blocking, per-sensor frame loops should be run on their own thread. This allows sensors to be processed at their own framerate.

3.2. Sensors and Sensor Frames

The `IResearchModeSensor` interface introduced in the previous section abstracts the Research Mode sensors. It provides methods and properties common to all sensors: `OpenStream`, `CloseStream`, `GetFriendlyName`, `GetSensorType`, `GetNextBuffer`. Sensors can be of the following types:

```

enum ResearchModeSensorType
{
    LEFT_FRONT,
    LEFT_LEFT,
    RIGHT_FRONT,
    RIGHT_RIGHT,
    DEPTH_AHAT,
    DEPTH_LONG_THROW,
    IMU_ACCEL,
    IMU_GYRO,
    IMU_MAG
};

```

where the first six types are *camera* sensors, and the remaining three are *IMU* sensors. Camera and IMU sensors differ in the methods they expose: for instance, camera sensors expose methods for projecting 3D points in camera space to 2D points in image space (see Sec. 3.3), while IMU sen-

sors do not. Sensor specializations are obtained by calling `QueryInterface`:

```

IResearchModeSensor* pSensor;
// ... initialize pSensor
IResearchModeCameraSensor* pCameraSensor = nullptr;
HRESULT hr = pSensor->QueryInterface(
    IID_PPV_ARGS(&pCameraSensor));

```

An analogous distinction exists between the `IResearchModeSensorFrame` interface and per-sensor specializations. Once a sensor is in streaming mode, sensor frames are retrieved with `IResearchModeSensor::GetNextBuffer`. All sensor frames have a common `IResearchModeSensorFrame` interface that returns frame information common to all types of frames: timestamps and sample size in bytes. As above, camera and IMU frame specializations are obtained by calling `QueryInterface`:

```

IResearchModeSensorFrame* pSensorFrame;
// ... Obtain pSensorFrame via GetNextBuffer

IResearchModeSensorVLCFrame* pVLCFrame = nullptr;
HRESULT hr = pSensorFrame->QueryInterface(
    IID_PPV_ARGS(&pVLCFrame));

```

Camera and IMU frames can be used to access the following information:

Camera frames provide getters for resolution, exposure, gain. In addition,

- VLC frames return grayscale buffers;
- Long Throw depth frames return a depth buffer, a sigma buffer and an active brightness buffer;
- AHAT depth frames return a depth buffer and an active brightness buffer.

The active brightness buffer returns a so-called IR reading. The value of pixels in the clean IR reading is proportional to the amount of light returned from the scene. The image looks similar to a regular IR image.

The sigma buffer for Long Throw is used to invalidate unreliable depth based on the invalidation mask computed by the depth algorithm. For AHAT, for efficiency purposes, the invalidation code is embedded in the depth channel itself.

The following code exemplifies how to access Long Throw depth:

```

void ProcessFrame(IResearchModeSensor *pSensor,
    IResearchModeSensorFrame* pSensorFrame,
    int bufferCount)
{
    ResearchModeSensorResolution resolution;
    ResearchModeSensorTimestamp timestamp;
    IResearchModeSensorDepthFrame *pDepthFrame = nullptr;
    UINT16 *pAbImage = nullptr;
    UINT16 *pDepth = nullptr;

```

```

BYTE *pSigma = nullptr;
// Invalidation mask for Long Throw
const USHORT mask = 0x80;

HRESULT hr = S_OK;
size_t outBufferCount;

pSensorFrame->GetResolution(&resolution);
pSensorFrame->GetTimeStamp(&timestamp);

hr = pSensorFrame->QueryInterface(
    IID_PPV_ARGS(&pDepthFrame));

if (SUCCEEDED(hr))
{
    hr = pDepthFrame->GetSigmaBuffer(&pSigma,
        &outBufferCount);
    // Process the buffer...
}
if (SUCCEEDED(hr))
{
    // Extract depth buffer
    hr = pDepthFrame->GetBuffer(&pDepth,
        &outBufferCount);
    // Validate depth
    for (size_t i = 0; i < outBufferCount; ++i)
    {
        // the most significant bit of pSigma[i]
        // tells if the pixel is valid
        bool isValid = (pSigma[i] & mask) > 0;
        if (isValid)
        {
            pDepth[i] = 0;
        }
    }
    // Process the buffer...
}
if (SUCCEEDED(hr))
{
    // Extract active brightness buffer
    hr = pDepthFrame->GetAbDepthBuffer(&pAbImage,
        &outBufferCount);
    // Process the buffer...
}
if (pDepthFrame)
    pDepthFrame->Release();
}

```

IMU frames store batches of sensor samples. Accelerometer frames store accelerometer measurements and temperature; gyroscope frames store gyroscope measurements and temperature; magnetometer frames store magnetometer measurements.

3.3. Sensor Coordinate Frames

All sensors are positioned in a device-defined coordinate frame, which is defined as the *rigNode*. Each sensor returns its transform to the *rigNode* (device origin) expressed as an extrinsics rigid body transform (rotation and translation). On HoloLens 2, the device origin corresponds to the Left Front Visible Light Camera; therefore, the transformation returned by this sensor corresponds to the identity transformation. Figure 5 shows camera coordinate frames relative to the *rigNode*.

The extrinsics transform for a camera sensor can be retrieved as follows:

```

IResearchModeCameraSensor *pCameraSensor;

```

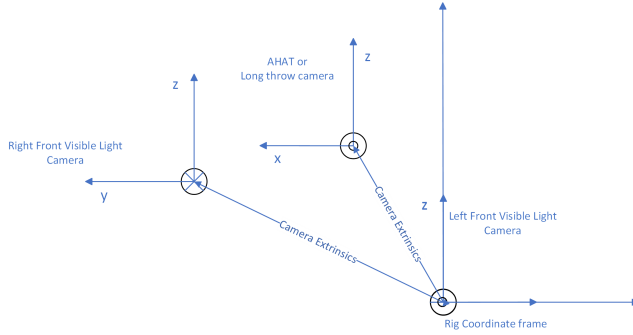


Figure 5. Camera coordinate frames relative to the *rigNode* (device origin – Left Front Visible Light camera).

```
DirectX::XMFLOAT4X4 cameraPose;
// Initialize camera sensor...
// Get matrix of extrinsics wrt the rigNode
pCameraSensor->GetCameraExtrinsicsMatrix(&cameraPose);
```

In order to locate sensors with respect to other coordinate systems (e.g. “in the world”), we use the HoloLens Perception APIs [7]. These APIs leverage the HoloLens head tracker to locate the device in a given coordinate system (for example, defined by the location of the user at app launch time). In order to locate the *rigNode*, we need to identify it with a GUID. This GUID is provided by the `IResearchModeSensorDevicePerception` interface:

```
using namespace winrt::Windows::Perception::Spatial;
using namespace winrt::Windows::Perception::Spatial::Preview;

HRESULT hr = S_OK;
SpatialLocator locator;
IResearchModeSensorDevicePerception*
    pSensorDevicePerception;
GUID guid;
hr = m_pSensorDevice->QueryInterface(
    IID_PPV_ARGS(&pSensorDevicePerception));
if (SUCCEEDED(hr))
{
    hr = pSensorDevicePerception->GetRigNodeId(&guid);
    locator =
        SpatialGraphInteropPreview::CreateLocatorForNode(guid);
}
// further processing, define anotherCoordSystem...
auto location = locator.TryLocateAtTimestamp(timestamp,
    anotherCoordSystem);
```

Note that camera sensors do not directly expose intrinsic parameters. Instead, they expose `MapImagePointToCameraUnitPlane` and `MapCameraSpaceToImagePoint` methods to convert 3D coordinates in the camera reference frame into 2D image coordinates, and vice versa.

4. The HoloLens2ForCV Repository

The HoloLens2ForCV repository [3] provides a few sample UWP (Universal Windows Platform) apps showing how to access and process Research Mode input streams:

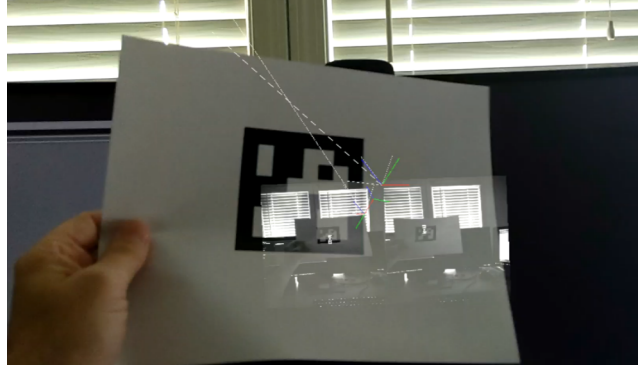


Figure 6. *CameraWithCVAndCalibration* app: We use OpenCV to detect arUco markers in the scene and triangulate them.

for example, how to visualize streams live on device, or how to record them and correlate them in space and time with hand and eye tracking. Section 4.1 provides a brief overview of the apps collected in the repository. Section 4.2 shows how, leveraging Research Mode and adding just a few lines of code, one can prototype computer vision algorithms like TSDF Volume Integration and SLAM.

4.1. Sample Apps

In its initial version, the HoloLens2ForCV repository collects four UWP apps. We plan to expand this initial set over time, and welcome contributions from the community. In particular, we provide sample code to a) visualize input streams live on device and process them using OpenCV [5] and b) record streams to disk for offline postprocessing.

Sensor Visualization and Processing: The *SensorVisualization* app allows one to visualize depth, VLC (and optionally IMU) input streams live on device (see Fig. 7). This can be particularly useful to test, for example, sensor performance in different environments. With the *CameraWithCVAndCalibration* app, we show how to process these streams on device: we use OpenCV to detect arUco markers in the two frontal VLC cameras, and triangulate the detections (Fig. 6). The *CalibrationVisualization* app can be used to visualize VLC and depth coordinate frames.

Stream Recorder: The *StreamRecorder* app allows one to capture simultaneously Research Mode streams (depth and VLC), the RGB stream from the HoloLens PV (PhotoVideo) camera, head, hand and eye tracking. The application provides a simple user interface to start and stop capture (Fig. 8, top). Streams are captured and stored on device, and can be downloaded and postprocessed by using a set of python scripts. Namely, we provide scripts to correlate Research Mode, PV, head, hand and eye tracking streams in



Figure 7. Research Mode sensor streams visualization on device: here we show the two frontal VLC cameras and Long Throw Depth.

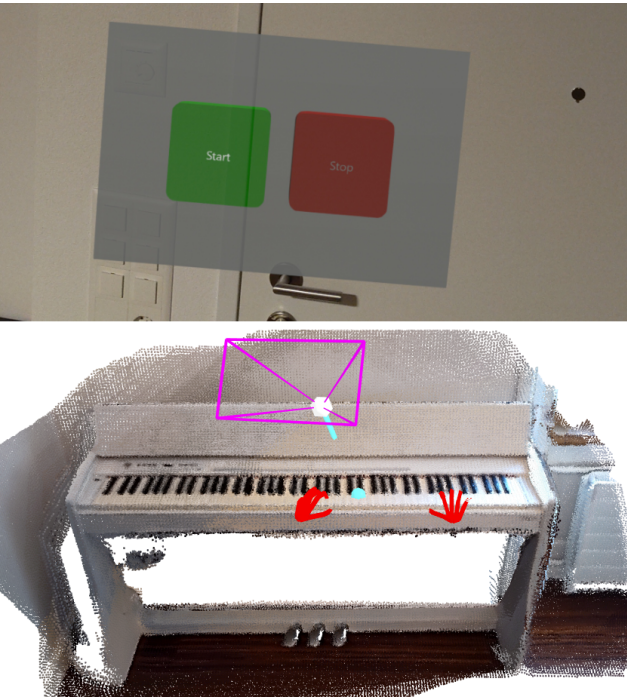


Figure 8. The *StreamRecorder* user interface (top). Processed outputs: Long Throw Depth frames (piano reconstruction), hand tracker output (red), eye gaze tracker output (blue sphere).

time and put them in a common coordinate frame. This allows us to reconstruct the scene in 3D and analyze how the user interacts with it by looking at their eye gaze, head and hand pose (Fig. 8, bottom). It is possible to add also audio capture, accessible via Windows Media APIs [6], though this functionality is not implemented at the time of writing.

It is worth noting that the *StreamRecorder* app relies on another library, *Cannon* (available on [3]). *Cannon* is a collection of wrappers and utility code for building native mixed reality apps using C++, Direct3D and Windows Per-

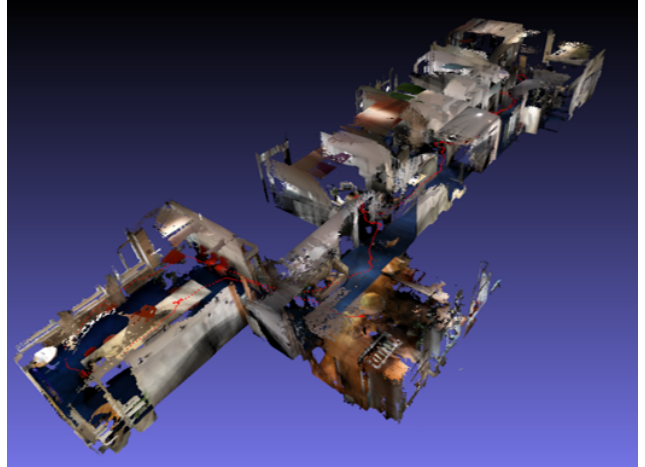


Figure 9. TSDF Volume Integration [9] result obtained from HoloLens depth and RGB input streams.

ception APIs. It can be used as-is outside Research Mode for fast and easy native development.

4.2. Computer Vision with Research Mode

By leveraging Research Mode, one can use the HoloLens as a powerful tool for computer vision and robotics research. To showcase how easily this can be done, we provide two examples: TSDF Volume Integration and SLAM.

TSDF Volume Integration: We provide a script which runs TSDF Volume Integration using as input Long throw depth frames, RGB frames and head poses captured with the *StreamRecorder* app. We simply combine the app output with an off-the-shelf library [9]. Figure 9 shows the reconstruction obtained by running the algorithm on almost 700 frames (head trajectory shown in red).

SLAM: We feed Long Throw depth frames and RGB frames as input to a recently proposed SLAM method (BAD SLAM [8]). The camera poses estimated by the SLAM algorithm can be compared and evaluated against the head poses returned by the HoloLens.

5. Conclusion

We presented Research Mode for HoloLens 2, an API and a set of functionalities enabling access to the raw sensor streams on device. Together with Research Mode, we release a public repository collecting sample apps and examples showcasing how to leverage Research Mode to build computer vision applications based on HoloLens. With these tools, we hope to facilitate further research in the fields of computer vision and robotics, and encourage contributions from the research community.

References

- [1] Computer Vision Applications for Mixed Reality Headsets Workshop at CVPR 2019. <https://docs.microsoft.com/en-us/windows/mixed-reality/cvpr-2019>. [Online; accessed August 2020].
- [2] Microsoft HoloLens2. <https://www.microsoft.com/en-us/hololens/hardware>. [Online; accessed August 2020].
- [3] Microsoft HoloLens2ForCV GitHub repository. <https://github.com/microsoft/HoloLens2ForCV>. [Online; accessed August 2020].
- [4] Microsoft HoloLensForCV GitHub repository. <https://github.com/microsoft/HoloLensForCV>. [Online; accessed August 2020].
- [5] OpenCV library. <https://opencv.org/>. [Online; accessed August 2020].
- [6] Windows Media APIs. <https://docs.microsoft.com/en-us/uwp/api/windows.media.capture.frames?view=winrt-19041>. [Online; accessed August 2020].
- [7] Windows Perception APIs. <https://docs.microsoft.com/en-us/uwp/api/windows.perception.spatial?view=winrt-19041>. [Online; accessed August 2020].
- [8] Thomas Schöps, Torsten Sattler, and Marc Pollefeys. BAD SLAM: Bundle adjusted direct RGB-D SLAM. In *IEEE CVPR*, 2019.
- [9] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. *arXiv:1801.09847*, 2018.