# Generating Benchmark Datasets for Computer Vision

**Daniel Scharstein**
Charles A. Dana Professor of Computer Science

**Toby Weed**
Middlebury class of 2022

**Motivation**
Computer vision is one of the most active research areas in computer science: think self-driving cars, facial and object recognition, image classification, and beyond.

In order to facilitate the development of better computer vision algorithms by hundreds of different teams across the world, it is necessary to know what "better" means.
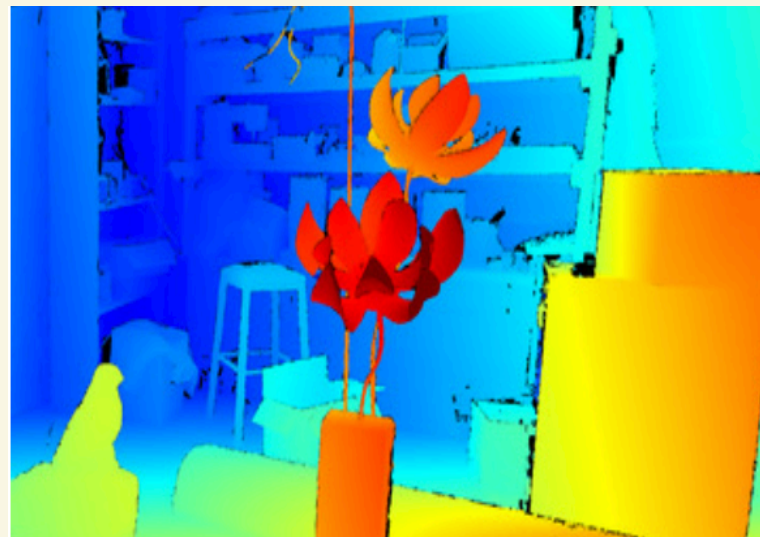
**The purpose of our research is to produce datasets which will help set the bar for computer vision algorithm performance.**

**Methods:** Researchers will run their algorithms to extract depth information from our scenes, and we'll compare their results with our own. Thus, a large part of our work is **extracting highly accurate 3D scans from pictures.**
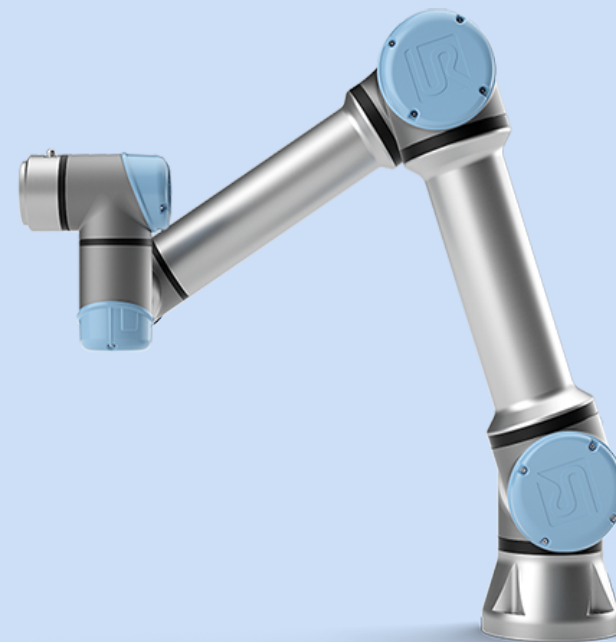
**Ground truth & Structured Lighting**

To get a highly accurate ground truth depth map of a scene, we project unique pixel codes on it and compare their motion.

**Robotic Viewpoint Control**

We use a UR5 industrial robot arm for repeatability and precision.

**Summer Progress & Future Work**
This is the fourth summer of research on this dataset, building on code from past interns. Here's what we worked on:

- **Replacing C++ calibration code** with a more flexible module, including a live visualization GUI.
- **Refining Swift, C++, & Python control servers** to expedite future maintenance, handover, and scene capture.
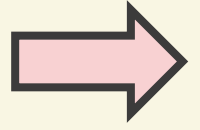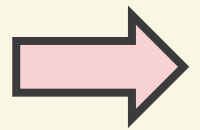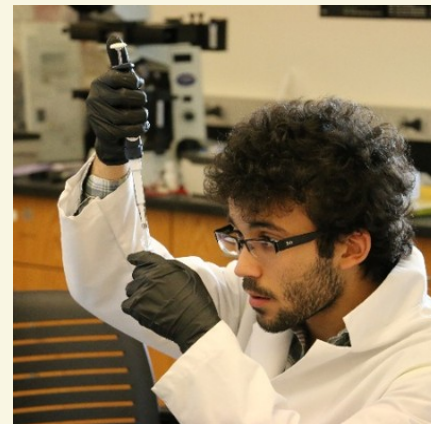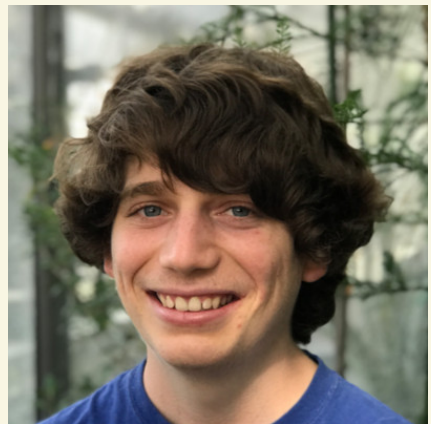- **Reimagining the end-goal** to best serve rapidly evolving evaluation needs.

We're hoping to continue moving forward during the academic year by integrating **realistic motion tracking** and **capturing production datasets**.

## Background & Future

This project has been worked on by numerous interns in the past:
- Guanghan Pan '21, Tiansheng Sun '20, me (summer '19)
- Nicholas Mosier '20, Tommaso Monaco '20, Roger Dai, '20 (summer '18)
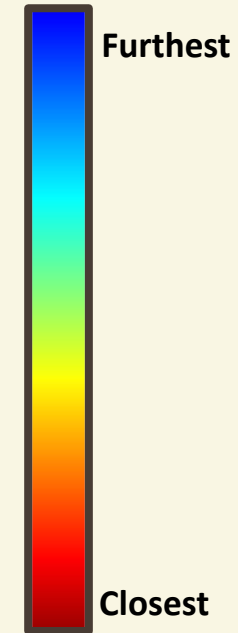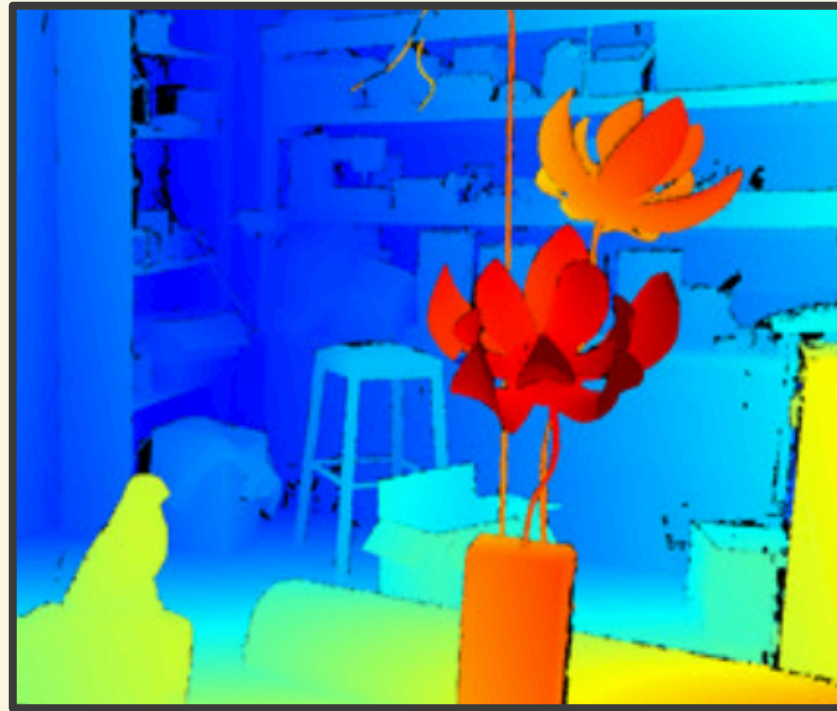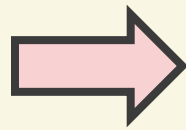- Nicholas Mosier'20 (summer '17)

I am continuing work this fall, but I expect (Professor Scharstein can confirm or deny) that more interns will be needed in summer '21).

.

# Motivation

Progress in stereo vision depends on the availability of precise methods for evaluating the accuracy of different algorithms.

This project will fill that void by providing diverse datasets with subpixel accurate ground truth depth maps against which researchers will be able to compare the results of their algorithms.
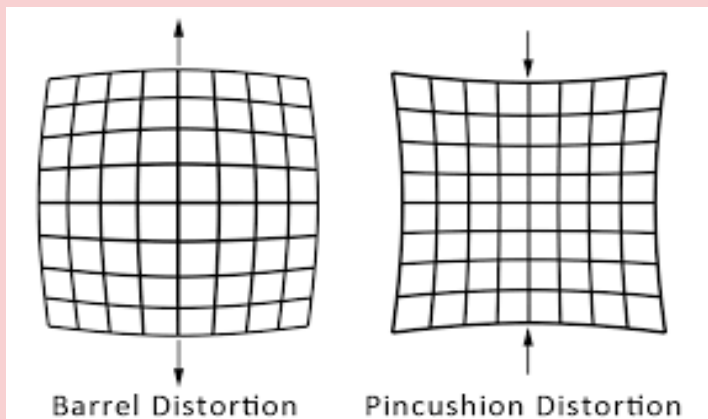


**So how do we get from here…**

**…to here?**

Furthest

Closest

## Example Use Cases for Mobile-Based Depth Perception

- **3D scanning of objects or faces:** e-commerce, Snapchat filters, security cameras.

- **Photo enhancement**: portrait mode, automatic editing effects.

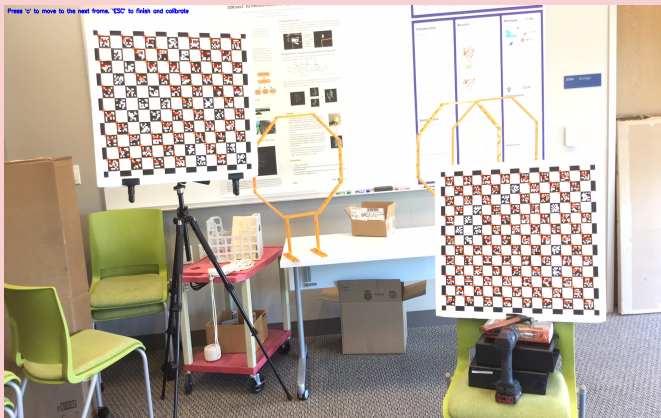- **AR/VR**: VR headsets, gaming, augmented video streaming.
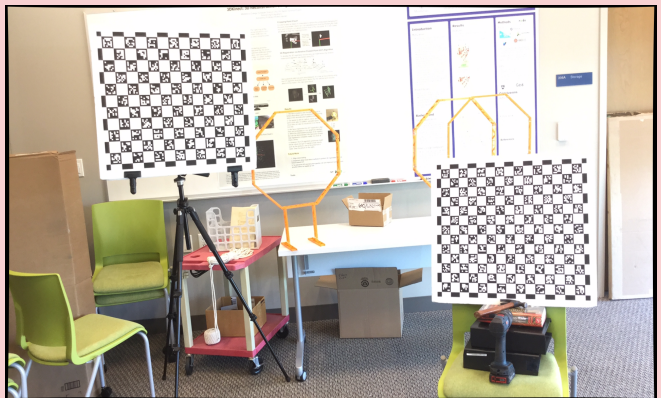
# Step One: Calibration



Lens Distortion

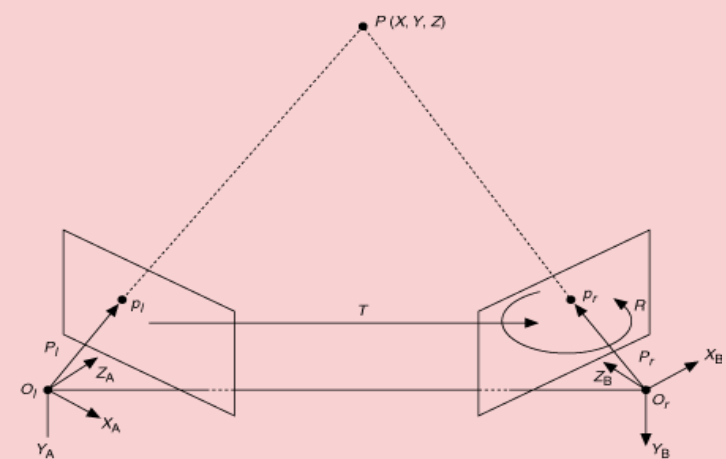Barrel Distortion     Pincushion Distortion

**Solved by** →

**Intrinsic Calibration**



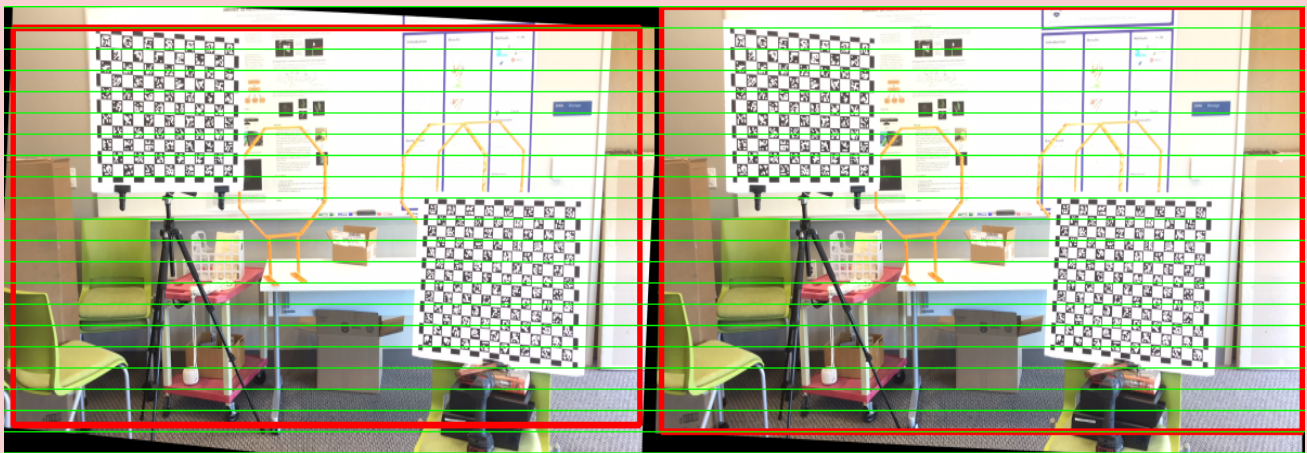**Marker Detection**



**Undistortion**



Viewpoint Disparities

**Solved by** →

**Extrinsic Calibration**



**Rectification**

Every camera lens causes some distortion, which we must account for if we seek to construct accurate 3D models. Additionally, we want to project images from each viewpoint onto a common plane to make finding their view disparities easier, which requires finding their locations & orientations relative to each other.

# Step Two: Structured Lighting & Disparity Matching

## Stereo Vision
Stereo vision is the extraction of 3D information from multiple 2D images taken from different locations. Our eyes do it and so can computers.
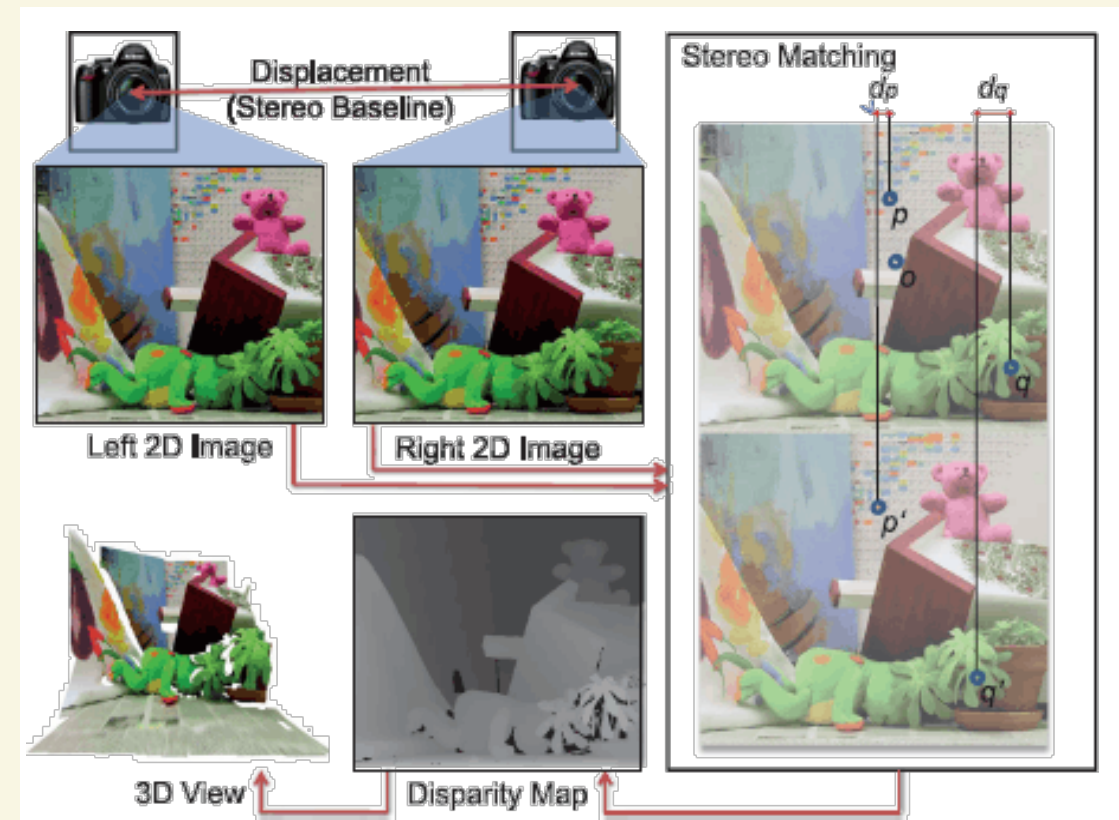
## Structured Lighting
The special sauce. By projecting structured lighting displays onto a scene, we can assign a **unique code to each projector pixel**, compare the movement of pixels between different viewpoints, and thus determine the proximity of a surface.

**More movement means the surface is closer, and less means it's farther.**



Projecting structured lighting onto a scene to obtain unique pixel codes for each point on a surface. We project 20 patterns, 10 with vertical stripes and 10 horizontal, to get a unique code for each projector pixel.



Comparing 2D codes of points on image to extract depth information.

# Step Three: Image Processing

To get ground truth depth maps of the highest possible quality, we use an image processing pipeline with 6 steps:

### Decode images

Decode the structured lighting images to get a pixel code for each point in the scene. The projected patterns act as a binary code for the X and Y coordinates.

### Rectify images

Using the extrinsic calibration parameters generated from calibration, rectify images from adjacent camera views so that they appear to lie on the same plane.

### Disparity match images

Merge the X and Y decoded images to get a unique pixel code for each point. Then compare the pixel codes from two views to get a disparity map.

### Merge disparity maps

Merge disparity maps produced by different projectors to fill in as many depth values as possible as accurately as possible.
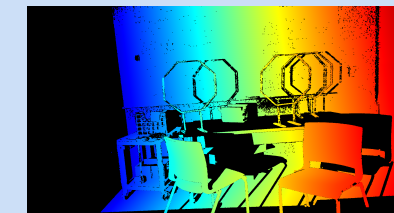
### Reproject merged disparity maps

The merged disparities are used to self-calibrate the projector positions. Once the projector relationships are known, half-occluded regions can be filled in with disparity values.
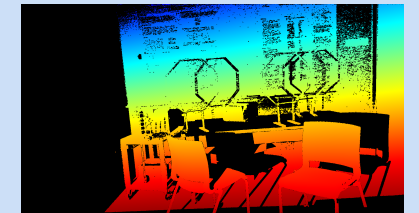
### Merge2

Put together original disparities, merged disparities, and reprojected disparities to fill in all possible values.
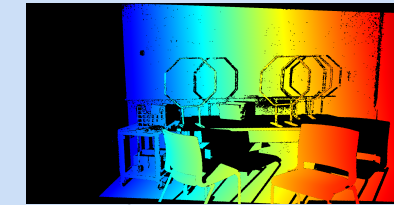
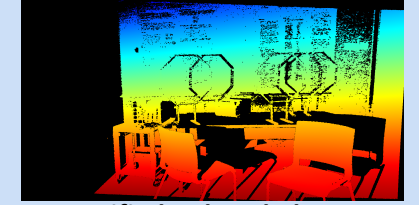These steps are all implemented using C++ in MobileLighting Mac.
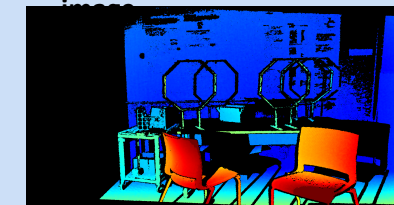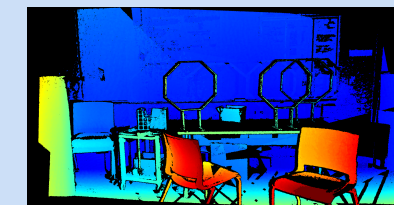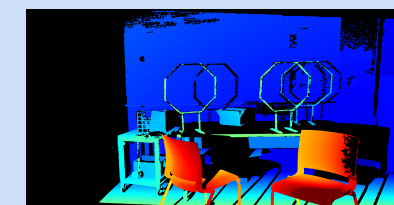


X-decoded image

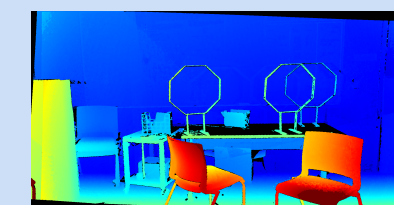Y-decoded image

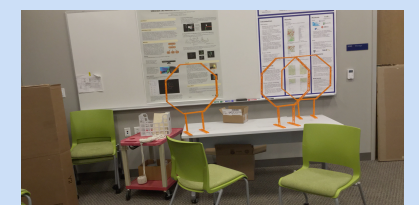Rectified X-decoded image

Rectified Y-decoded image

Initial disparity map

Merged disparity map

Reprojected disparity map

Final disparity map

Ambient image of the scene

**The image processing pipeline**

# Contributions this summer

## Marker Detection Visualization & Calibration Refactor

```
Hit Enter to begin taking photos, or q then enter to quit.

Taking a photo
-- CameraServiceBrowser: Sending packet #4530292560 with timestamp 47:7:574918031
-- PhotoReceiver: Read data with tag 1
-- PhotoReceiver: packetDataLength: 1310155
-- PhotoReceiver: Read data with tag 2
-- PhotoReceiver: Successfully saved photo data to file file:///Volumes/My%20Passport//sandbox/orig/calibration/intrinsics/IMG0.JPG.
-- PhotoReceiver: Received calibration image.
Tracking ChArUco markers from image
Reading board 0 from file /Volumes/My Passport//sandbox/settings/boards/board0.yml
Reading from file /Volumes/My Passport//sandbox/settings/boards/board0.yml
Reading board 1 from file /Volumes/My Passport//sandbox/settings/boards/board1.yml
Reading from file /Volumes/My Passport//sandbox/settings/boards/board1.yml

Reading image from file /Volumes/My Passport//sandbox/orig/calibration/intrinsics/IMG0.JPG
Detecting ArUco markers
Interpolating chessboard corners from board 0 based on detected ArUco markers
Interpolating chessboard corners from board 1 based on detected ArUco markers
Drawing detected marker indicators
Drawing chessboard corners

With image display window open, press any key to continue, r to retake, or q to quit.
```
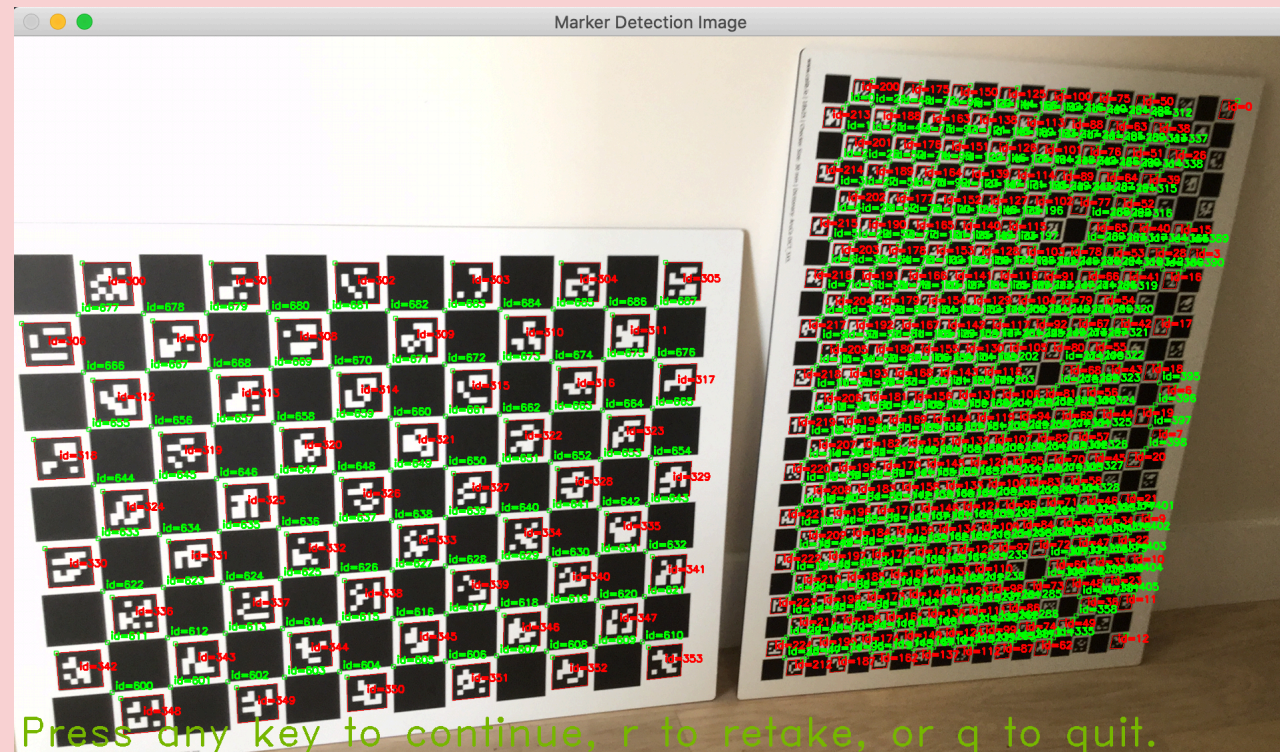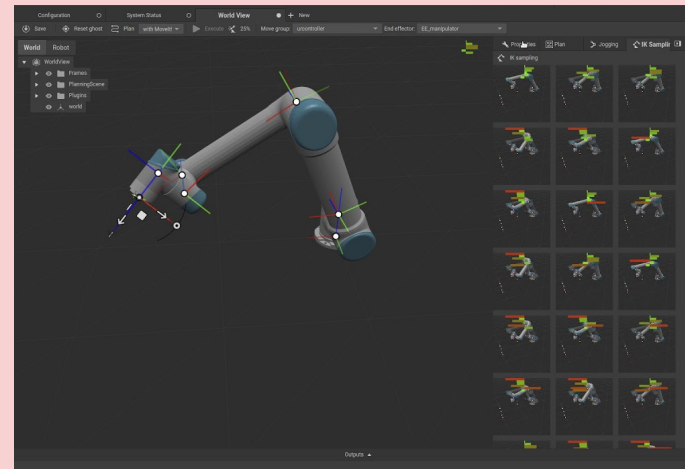


## Universal Robot State Tracking



MLRobotControl server (Python)

Trackfile for ML_Mac system coordination server (Swift & C++)

## Full-System Usability, Modularity, Style Upgrades

- Integrated processing mode for detachability
- Improved usage messages
- Fixed bugs
- Improved package management and portability
- Standardized input and output files

## Challenges

- Large, unwieldy codebase
- Tenacious iPhone and Mac communication issues
- Inability of robot to capture smooth video
- Mixed languages, protocols, and programming paradigms

# References and Credits

**Learn More**
Middlebury Stereo Homepage:
http://vision.middlebury.edu/stereo/
Auto-generated web pages displaying some samples from datasets we took last summer:
shorturl.at/gEMRZ
shorturl.at/dBGMQ
Project GitHub (contains tons of information and links to other relevant things):
https://github.com/tobyweed/MobileLighting

**Acknowledging to prior Middlebury research assistants, without whom the system wouldn't exist:**
- Guanghan Pan '21
- Tiansheng Sun '20
- Nicholas Mosier '20
- Tommaso Monaco '20
- Roger Dai, '20

**Thank you to Professor Scharstein, the CS department, and the NSF!**