

# Mizzou INformation and Data FUsion Lab (MINDFUL)

**Title:** Frame Selection Strategies for Real-Time Structure-from-Motion from an Aerial Platform

**Authors:** Andrew R. Buck, Jack Akers, Derek T. Anderson, James M. Keller,  
Raub Camaioni, Matthew Deardorff, and Robert H. Luke III

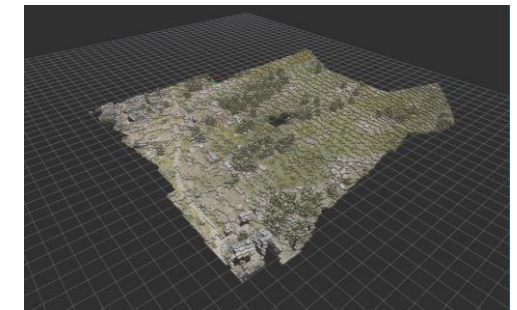


**University of Missouri**



# Introduction

- **What's our problem?**
  - Real-time 3D mapping from a single camera on a micro-UAV
- **Challenges?**
  - Single camera
  - Can solve with SfM on frame pairs
  - Don't control where drone goes (manual control)
  - How to select frames to obtain the best reconstruction?
- **Today**
  - SIM environment and study







# Random Movement Dataset

UE native  
quality



Reconstructed  
3D voxel space  
(UFOmap)



AirSim  
quality



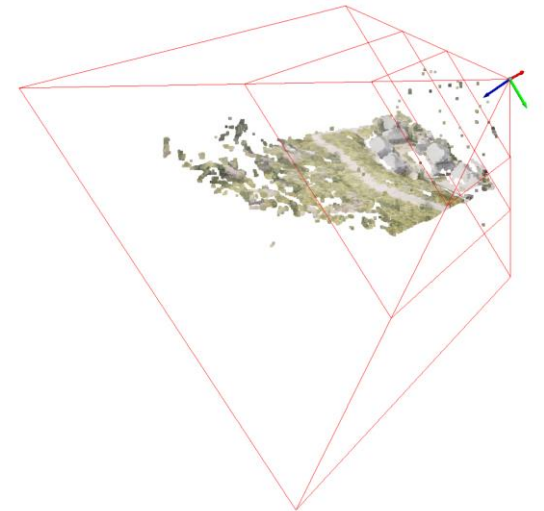
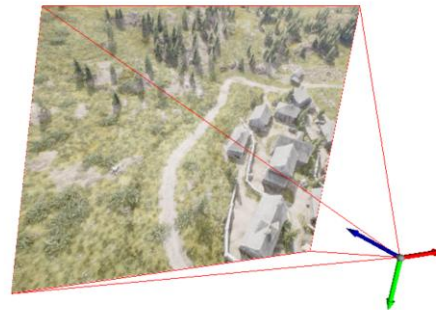
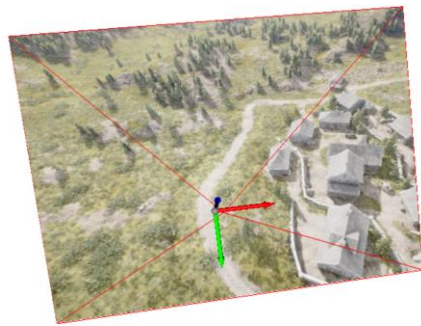
Collected about 1000  
frames of aligned  
RGB/Depth imagery

Drone moves to  
random poses  
(position and look  
direction)



# EpiDepth

- A moving camera on a UAV provides a stream of images with known poses (thanks to onboard GPS/IMU).
- For a given frame pair, we can align the images and perform stereo matching to estimate depth.

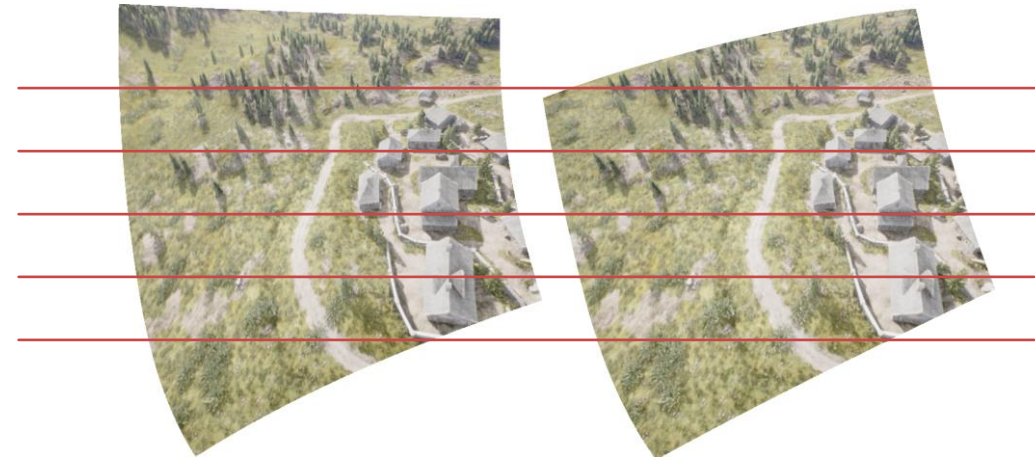
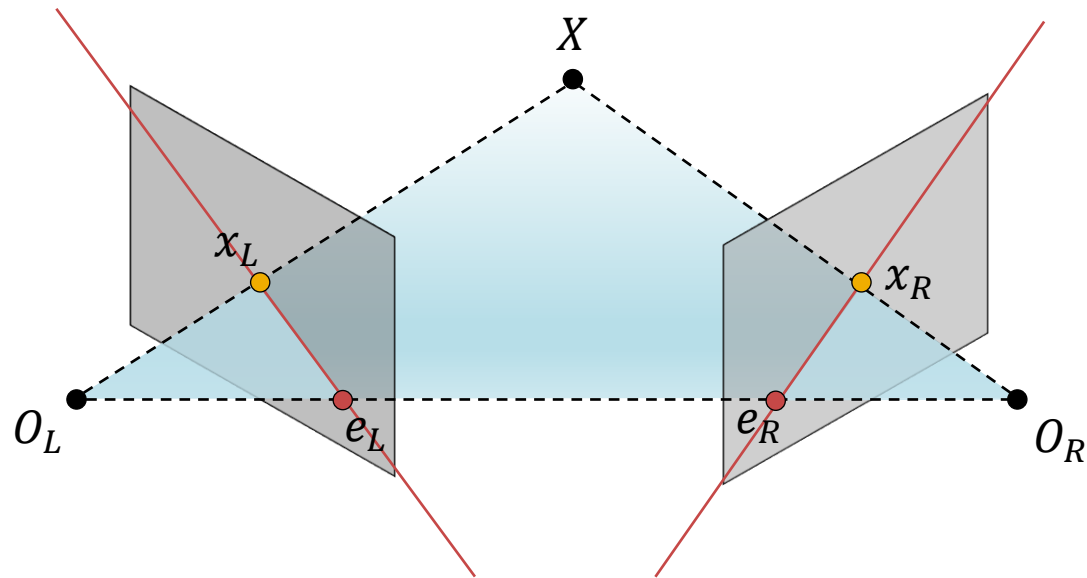






# Epipolar Warping

- The epipolar geometry of two camera views defines how to warp the images.
- Feature pairs are aligned on the same row and the pixel disparity is used to estimate depth.





# Warping Effects

- The relative pose between images has a big impact on how much warping is required.
- Generally, areas around the epipole are hard to match.

Nadir



Climb



Forward 45°



Straight Ahead







# EpiDepth Matching Examples

- Good cases for frame pair matching...

Error Metrics	Camera Extrinsic	Frame 1 Warped	Frame 2 Warped	Frame 1	Frame 2	Ground Truth Depth	Predicted Depth	Depth Error
completeness: 0.548 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 1.121 rmse_log: 0.024 abs_rel: 0.017 sq_rel: 0.027	 EM: 0.997 BL: 2.360							
completeness: 0.565 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 0.787 rmse_log: 0.019 abs_rel: 0.013 sq_rel: 0.015	 EM: 0.953 BL: 3.017							
completeness: 0.488 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 1.861 rmse_log: 0.027 abs_rel: 0.019 sq_rel: 0.049	 EM: 0.707 BL: 3.072							
completeness: 0.309 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 1.558 rmse_log: 0.024 abs_rel: 0.016 sq_rel: 0.035	 EM: 0.555 BL: 2.898							





# EpiDepth Matching Examples

- Not so good cases for frame pair matching...

Error Metrics	Camera Extrinsic	Frame 1 Warped	Frame 2 Warped	Frame 1	Frame 2	Ground Truth Depth	Predicted Depth	Depth Error
completeness: 0.081 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 1.526 rmse_log: 0.022 abs_rel: 0.015 sq_rel: 0.031	 EM: 0.353 BL: 3.813							
completeness: 0.629 a1: 0.371 a2: 0.718 a3: 0.922 rmse: 16.427 rmse_log: 0.398 abs_rel: 0.291 sq_rel: 5.731	 EM: 0.847 BL: 0.048							
completeness: 0.057 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 1.219 rmse_log: 0.014 abs_rel: 0.011 sq_rel: 0.017	 EM: 0.977 BL: 14.249							
completeness: 0.000 a1: nan a2: nan a3: nan rmse: nan rmse_log: nan abs_rel: nan sq_rel: nan	 EM: 0.521 BL: 19.113							





# Simple Frame Picker (Version 1)

- Keeps a running candidate frame and tries to match new incoming frames to this one.
- Reject if frames are too close or if the rotation difference is too large.
- If a pair is found, yield it and keep the latest frame as the new frame to match to.
- Otherwise, if distance becomes too large, replace the candidate frame with the latest frame.

## Frame Picker v1

Define  $d_{\min}$ ,  $d_{\max}$ , and  $r_{\max}$

Initialize frame  $f_0$

For each new frame  $f_t$ :

// Get baseline distance

$d \leftarrow \text{DISTANCE}(f_0, f_t)$

// Check if extrinsics are acceptable

If  $d_{\min} < d < d_{\max}$ :

$r \leftarrow \text{ROTATION}(f_0, f_t)$

If  $r < r_{\max}$ :

yield  $(f_0, f_t)$

$f_0 \leftarrow f_t$

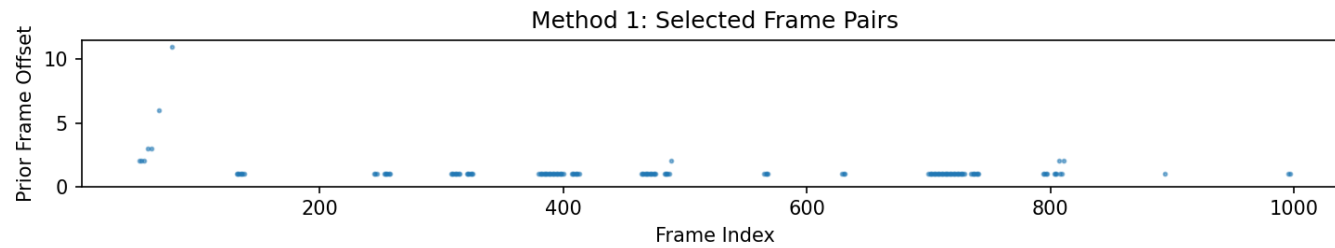
// Drop old frame

If  $d > d_{\max}$ :

$f_0 \leftarrow f_t$



# Method 1 (Simple) Results



Total frame pairs selected: **142**  
 Completeness:  $0.432 \pm 0.226$   
 RMSE-Log:  $0.035 \pm 0.010$

## Examples:

completeness: 0.589 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 1.473 rmse_log: 0.033 abs_rel: 0.024 sq_rel: 0.047	 EM: 1.000 BL: 1.064								
completeness: 0.640 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 1.658 rmse_log: 0.027 abs_rel: 0.020 sq_rel: 0.042	 EM: 0.979 BL: 2.019								
completeness: 0.679 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 1.880 rmse_log: 0.033 abs_rel: 0.025 sq_rel: 0.059	 EM: 0.999 BL: 1.007								
completeness: 0.200 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 2.263 rmse_log: 0.033 abs_rel: 0.024 sq_rel: 0.072	 EM: 0.483 BL: 1.651								



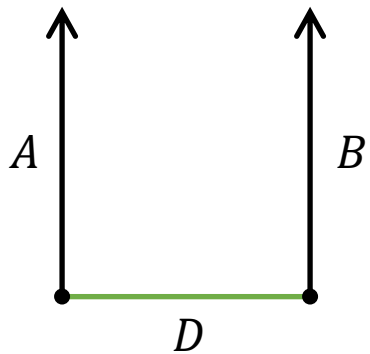


# Extrinsic Quality Metric

- We can define some heuristics to judge the quality of the two frame poses.
  - Let  $A$  and  $B$  be the look vectors of the two image frames
  - Let  $D$  be the displacement between the focal points of the two image frames

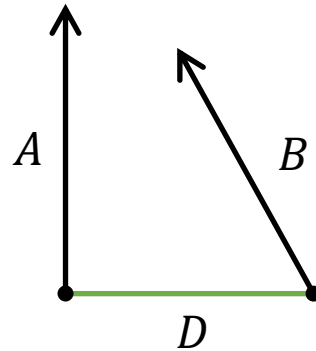
Heuristics:

- $\angle AB$  should be small
- $\angle AD$  and  $\angle BD$  should both be close to  $90^\circ$



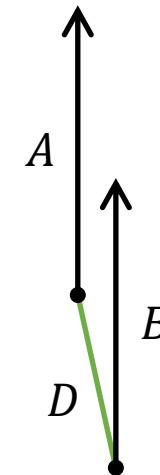
Excellent

Frames are separated perpendicular to the look direction and aligned



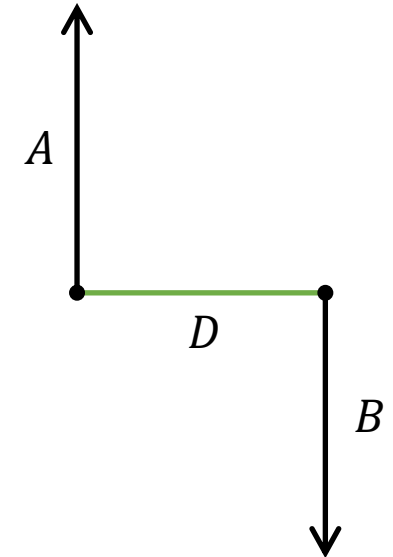
Good

Frames are separated and mostly aligned



Poor

Frames are aligned, but separated in the look direction



Bad

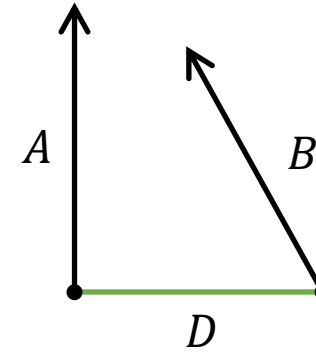
Frames are not aligned



# EQ Metric Function Crafting

- $\angle AB$  should be small (●)

$$S_{AB} = \cos(\angle AB)$$
$$H_{AB} = \begin{cases} 0, & S_{AB} < 0 \\ S_{AB}, & S_{AB} \geq 0 \end{cases}$$



- $\angle AD$  and  $\angle BD$  should both be close to  $90^\circ$  (●)

$$S_{AD} = \cos(\angle AD)$$

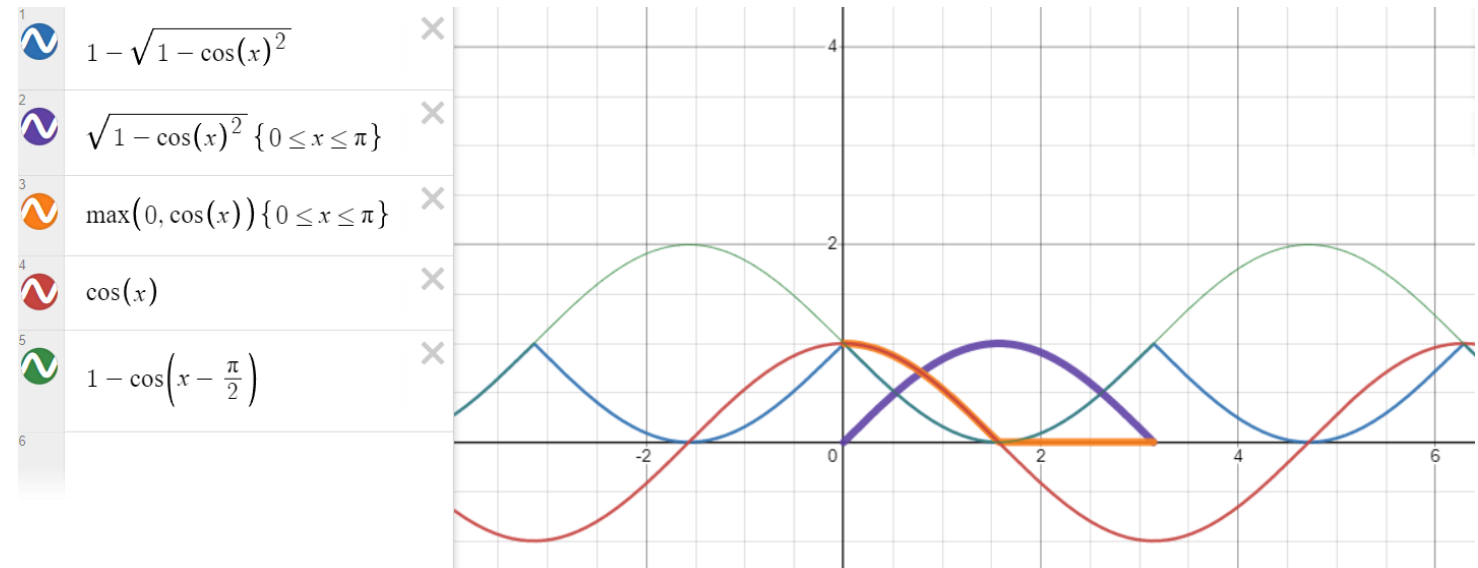
$$S_{BD} = \cos(\angle BD)$$

$$R_{AD} = \sqrt{1 - S_{AD}^2}$$

$$R_{BD} = \sqrt{1 - S_{BD}^2}$$

- Overall metric is the minimum of these,

$$Q_{ABD} = \min(H_{AB}, R_{AD}, R_{BD})$$







# Extrinsic Quality Examples

<u>Extrinsic Quality</u>	<u>Top</u>	<u>Front</u>	<u>Side</u>	<u>Frame 1</u>	<u>Frame 2</u>	
1.000						Nadir
1.000						Strafe
0.707						45° Climb
0.680						45° Forward
0.001						Straight Ahead
0.000						Nadir Climb



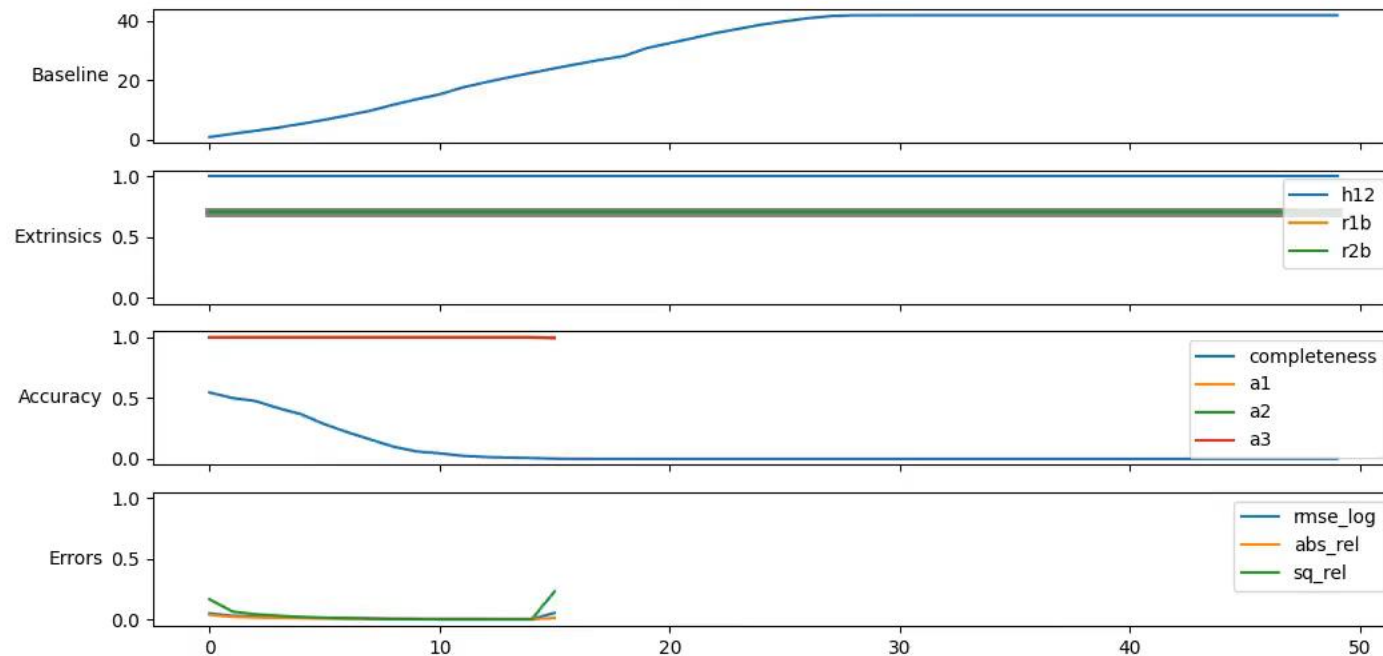
# Rolling Frame Buffer

As new frames arrive, add them to the buffer.

Compute baseline distance and extrinsic metric with each frame in the buffer.

Use these features to select the best frame to match with.

Buffer for Frame 50







# Heuristic Frame Picker (Version 2)

- Keep rolling buffer, past  $N$  frames.
- For each new frame, look back in the buffer for the most recent frame that has a minimum baseline distance.
- If this frame has an acceptable extrinsic metric when compared with the current frame, yield the frame pair.
- Otherwise, move on to the next incoming frame.

## Frame Picker v2

Define  $N$ ,  $d_{\min}$ ,  $q_{\min}$

Initialize rolling frame buffer  $B$

For each new frame  $f_t$ :

$B.insert(f_t)$

    For  $i = 1 \dots N$ :

$d \leftarrow \text{DISTANCE}(f_{t-i}, f_t)$

        If  $d < d_{\min}$ :

            continue

$q \leftarrow \text{EXTRINSIC\_QUALITY}(f_{t-i}, f_t)$

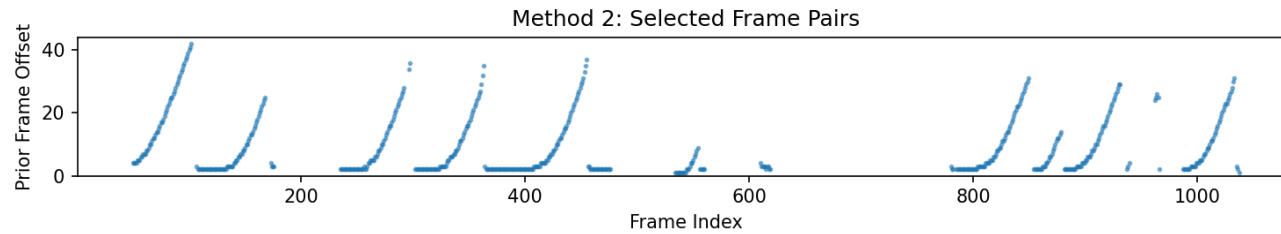
        If  $q \geq q_{\min}$ :

            yield  $(f_{t-i}, f_t)$

        break



# Method 2 (Heuristic) Results



Total frame pairs selected: **592**  
 Completeness:  $0.364 \pm 0.190$   
 RMSE-Log:  $0.030 \pm 0.138$

## Examples:

completeness: 0.281 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 2.411 rmse_log: 0.036 abs_rel: 0.026 sq_rel: 0.082	 EM: 0.900 BL: 4.106							
completeness: 0.074 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 2.702 rmse_log: 0.031 abs_rel: 0.022 sq_rel: 0.082	 EM: 0.958 BL: 3.212							
completeness: 0.543 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 0.639 rmse_log: 0.015 abs_rel: 0.010 sq_rel: 0.009	 EM: 0.957 BL: 4.559							
completeness: 0.678 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 0.860 rmse_log: 0.019 abs_rel: 0.012 sq_rel: 0.016	 EM: 0.976 BL: 3.646							

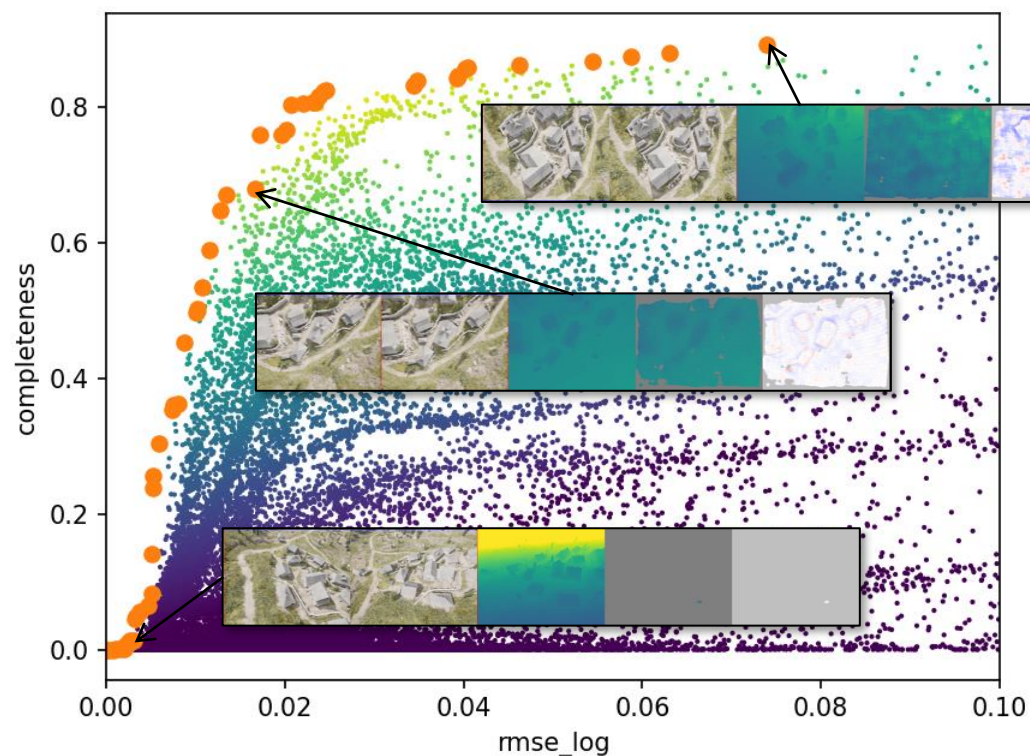




# Evaluating All Frame Pairs

- Since we have computed the EpiDepth prediction for all possible frame pairs, can we make use of this data?
- We want results to have low RMSE-log values and high completeness.
- Looking at examples from the Pareto front, we select appropriate scalarization weights.
- A neural net is trained on this dataset and can be used to predict if a pose pair will perform well.

RMSE-Log vs. Completeness of EpiDepth on All Frame Pairs





# Data-Driven Frame Picker (Version 3)

- Keep a rolling buffer of the past  $N$  frames.
- For each new frame, evaluate the predicted quality with all frames in the buffer.
- If the best scoring frame satisfies the minimum baseline and extrinsic quality thresholds, yield it.
- Otherwise, move on to the next incoming frame.

## Frame Picker v3

Define  $N, d_{\min}, q_{\min}, p_{\min}$

Initialize rolling frame buffer  $B$

For each new frame  $f_t$ :

$B.insert(f_t)$

$f_{\text{best}} \leftarrow \emptyset$

$p_{\text{best}} \leftarrow -\infty$

For  $i = 1 \dots N$ :

$d \leftarrow \text{DISTANCE}(f_{t-i}, f_t)$

$q \leftarrow \text{EXTRINSIC\_QUALITY}(f_{t-i}, f_t)$

If  $d < d_{\min}$  or  $q < q_{\min}$ :

continue

$p \leftarrow \text{PREDICT}(f_{t-i}, f_t)$

If  $p > p_{\text{best}}$  and  $p > p_{\min}$ :

$p_{\text{best}} \leftarrow p$

$f_{\text{best}} \leftarrow f_{t-i}$

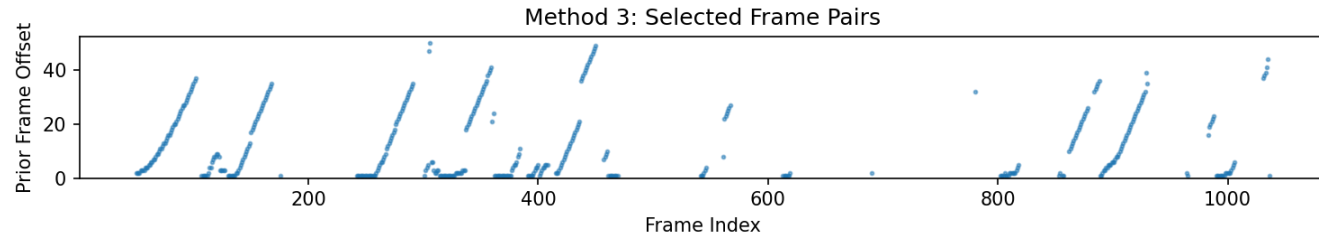
If  $f_{\text{best}} \neq \emptyset$ :

yield  $(f_{\text{best}}, f_t)$





# Method 3 (Data-Driven) Results



Total frame pairs selected: **453**  
 Completeness:  $0.322 \pm 0.249$   
 RMSE-Log:  $0.024 \pm 0.027$

## Examples:

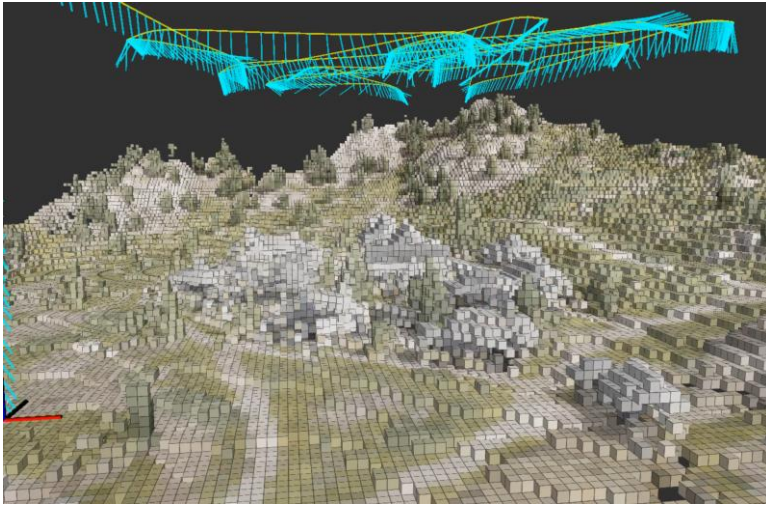
completeness: 0.050 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 1055.396 rmse_log: 0.118 abs_rel: 0.007 sq_rel: 17.009	 EM: 0.953 BL: 17.961						
completeness: 0.799 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 1.235 rmse_log: 0.027 abs_rel: 0.019 sq_rel: 0.033	 EM: 0.980 BL: 1.569						
completeness: 0.645 a1: 1.000 a2: 1.000 a3: 1.000 rmse: 1.818 rmse_log: 0.029 abs_rel: 0.022 sq_rel: 0.050	 EM: 0.978 BL: 1.536						
completeness: 0.213 a1: 0.997 a2: 1.000 a3: 1.000 rmse: 3.913 rmse_log: 0.058 abs_rel: 0.044 sq_rel: 0.209	 EM: 0.969 BL: 1.448						



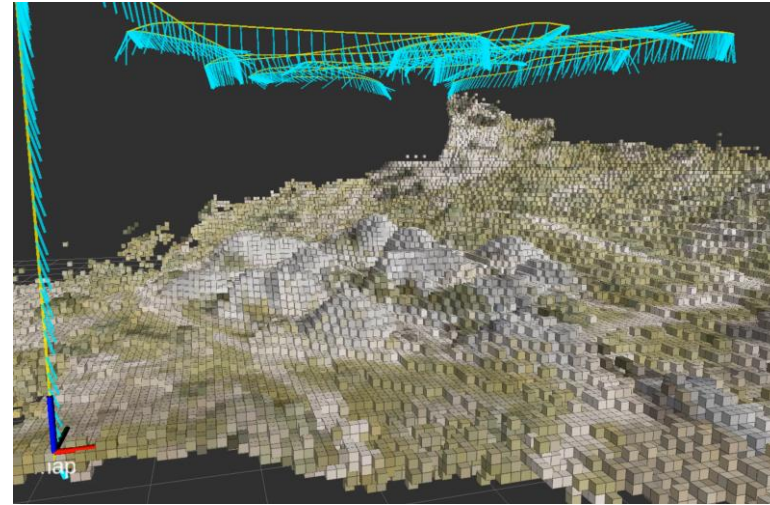


# Qualitative 3D Evaluation

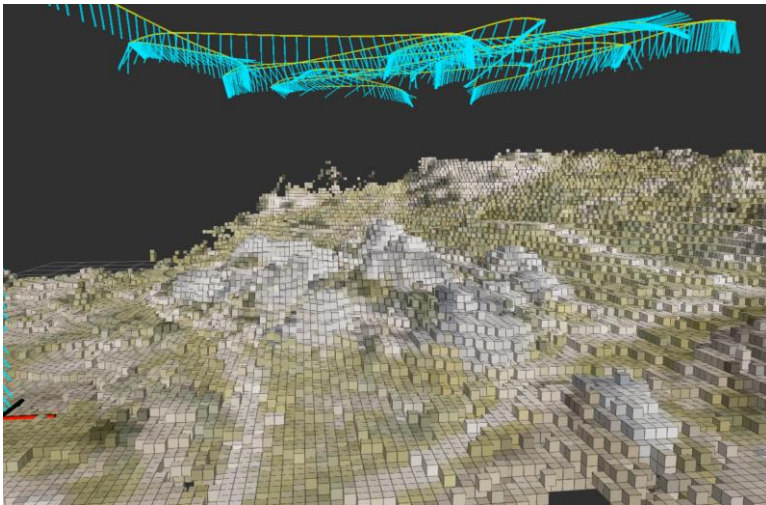
Ground Truth



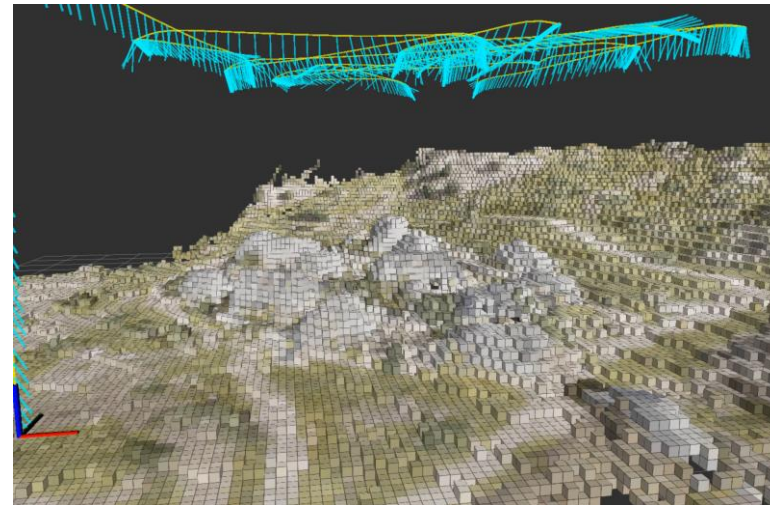
Method 1  
(Simple)



Method 2  
(Heuristic)



Method 3  
(Data-Driven)

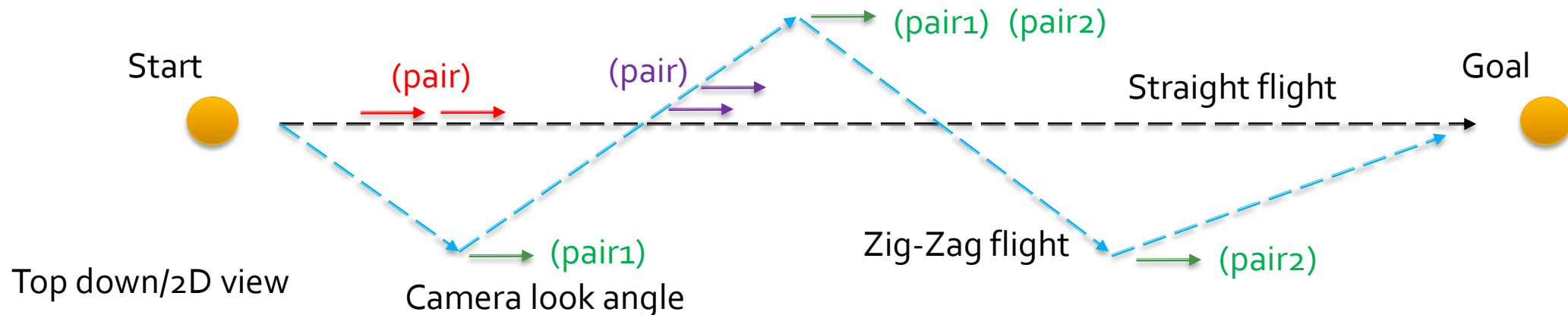






# Trust These SIM Results?

- Well, as real as SIM (and our setup at that) is real.
  - What are the important variables to SIM?
- 1. **Extrinsic error**
  - BIG real-world problem for sure! (position, pose, ...)
  - Can drastically impact UFOmap (our 3D aggregate structure), even when/if EpiDepth's good.
  - Does this mean “fewer projected points are more dangerous...”?
- 2. **Specific behaviors**
  - We ran a random movement sequence, wanted to remove any bias.
  - But did that wash out behavior specific benefits?
    - e.g., fly a zig-zag vs straight flight pattern, larger baselines, see further out





# Summary

- **SIM framework for SfM frame selection**
  - Real-time micro-UAV
  - Code base is nearly one-2-one with real platform
- **Explored: (1) naïve, (2) heuristic, (3) data-driven**
  - Best-2-worse → data-driven (3) then (2) then (1)
  - Quant: image space error, # of frames picked, % completeness
  - Qual: 3D voxel space reconstruction and 2D visualizations
- **Induced too easy/clean of a SIM setup?**
  - Expected bigger differences (like we have been seeing on real)
  - Perfect extrinsic data and random flight pattern
  - Goal was to generate a dataset that covers the distribution of all possible pose configurations and study in detail





# Next Steps

- **SIM environment**
  - Suspected we might need to include more factors
    - Add extrinsic error and specific flight patterns to SIM
- **SIM experiments**
  - Quantitative 3D voxel space metrics (from our SPIE paper)
    - How good at free space, occupied, time-varying, etc.
  - Decomposition by “variables” (from our MSS paper)
    - How good w.r.t. object ID type, breakdown by range, etc.
- **Frame selection**
  - Preliminary work
  - More in-depth analysis of these three real-time solutions
  - Improved data-driven solution
- **SIM vs real**
  - Real world confirmational experiments



# Questions?