

# Scalable Structure from Motion for Densely Sampled Videos

## Supplemental Material

B. Resch<sup>1,2</sup>   H. P. A. Lensch<sup>2,3</sup>   O. Wang<sup>1</sup>   M. Pollefeys<sup>3</sup>   A. Sorkine-Hornung<sup>1</sup>  
<sup>1</sup>Disney Research Zurich   <sup>2</sup>Tübingen University   <sup>3</sup>ETH Zurich

In this supplemental material, we present some more information about our test datasets. Additionally, we show results from reconstructions of datasets not used for evaluation in the paper: Stanford Light Fields Archive datasets, 5k video and a cooperative data acquisition scenario where the same room was captured with different cameras.

### 1. Paper dataset information

We used the following datasets in the paper for which we show some example images and illustrations of the reconstructions, as well as individual reconstruction timings (Figure 1).

**SHIP** The SHIP dataset consists of a video sequence with 4353 frames, recorded with a DSLR camera in 1080p. The camera was mounted on a slowly moving crane to simulate high framerates: the camera translation from frame to frame is about 1/4 of the other test scenes in the paper, so it corresponds to the same camera velocity at about 240 fps. The trajectory starts and stops at the same position which allows camera drift to be measured. Refer to Figure 2 for some example images and a 3D visualization of the reconstruction.

**OFFICE** The OFFICE dataset consists of a video sequence with 1333 frames, recorded with a GoPro Hero 3 camera in 1080p 60 fps wide angle mode. The trajectory starts and stops at the same position which allows camera drift to be measured. Refer to Figure 3 for some example images and a 3D visualization of the reconstruction.

**BENCH** The BENCH dataset consists of a video sequence with 1055 frames, recorded with a GoPro Hero 3 camera in 1080p 60 fps wide angle mode. The trajectory starts and stops at the same position which allows camera drift to be measured. Refer to Figure 4 for some example images and a 3D visualization of the reconstruction.

**STAIRWAY** The STAIRWAY dataset consists of a video sequence with 876 frames, recorded with a GoPro Hero 3 camera in 1080p 60 fps wide angle mode. The trajectory

starts and stops at the same position which allows camera drift to be measured. Refer to Figure 5 for some example images and a 3D visualization of the reconstruction.

**OFFICE 2** The OFFICE 2 dataset consists of a video sequence with 1180 frames, recorded with a GoPro Hero 3 camera in 1080p 60 fps wide angle mode. The trajectory starts and stops at the same position which allows camera drift to be measured. Refer to Figure 6 for some example images and a 3D visualization of the reconstruction.

**SWINGING** The SWINGING dataset consists of a video sequence with 1257 frames, recorded with a GoPro Hero 3 camera in 1080p 60 fps wide angle mode. The trajectory starts and stops at the same position which allows camera drift to be measured. Refer to Figure 7 for some example images and a 3D visualization of the reconstruction.

**2 LOOPS** The 2 LOOPS dataset consists of a video sequence with 1715 frames, recorded with a GoPro Hero 3 camera in 1080p 60 fps wide angle mode. The trajectory starts and stops at the same position which allows camera drift to be measured. The camera circles an object twice which allows to compare our algorithm with and without global anchor constraints between the two loops. Refer to Figure 8 for some example images and a 3D visualization of the reconstruction.

### 2. Pseudocode

For easier reimplemention, we give the pseudo code of our method in Algorithm 1.

### 3. Stanford lightfields

We reconstructed some of the Stanford Light Field Archive datasets by treating the individual 10 megapixel raw images in zigzag order as video stream. Each dataset contains an array of 17x17=289 camera poses. For timings refer to Figure 9. An exemplary reconstruction visualization of the lego truck scene can be found in Figure 10.

---

**Algorithm 1** High level pseudo code of the proposed method. This algorithm covers the main function "reconstruct" and a selected subset its called functions.

---

```

1: function RECONSTRUCT(images, costMatrix)
2:   windows  $\leftarrow$   $\emptyset$  ▷ windows only contains keyframes
3:   for all imageSeq  $\subseteq$  images do
4:     tracks  $\leftarrow$  TRACKING(imageSeq)
5:     startPoses  $\leftarrow$  INITWINDOWBA(imageSeq, tracks) ▷ See paper Section 3.2.1.
6:     windows  $\leftarrow$  windows  $\cup$  WINDOWBA(startPoses, imageSeq, tracks)
7:   end for
8:   windows  $\leftarrow$  windows  $\cup$  ANCHORHANDLING(windows, costMatrix)
9:   constraints  $\leftarrow$  EXTRACTCONSTRAINTS(windows) ▷ constraints only contains information about keyframes
10:  scene  $\leftarrow$  INITPOSESLINEAR(constraints) ▷ scene only contains keyframes poses here
11:  scene  $\leftarrow$  OPTIMIZE(scene)
12:  scene  $\leftarrow$  INTERPOLATENONKEYFRAMEPOSES(scene)
13:  scene  $\leftarrow$  OPTIMIZE(scene)
14: end function
15:
16: function TRACKING(images) ▷ See paper Section 3.1.
17:  allTracks  $\leftarrow$   $\emptyset$ 
18:  activeTracks  $\leftarrow$   $\emptyset$ 
19:  for image  $\in$  images do
20:    for all track  $\in$  activeTracks do
21:      feature  $\leftarrow$  FINDFEATURE(image, track.previousImagePattern)
22:      feature  $\leftarrow$  REFINEFEATURE(feature, image, track.firstImagePattern)
23:      if feature  $\neq$  null then
24:        EXTENDTRACK(track, feature)
25:      else
26:        activeTracks  $\leftarrow$  activeTracks  $\setminus$  {track}
27:      end if
28:    end for
29:    newTracks  $\leftarrow$  FINDNEWKEYPOINTS(image)
30:    newTracks  $\leftarrow$  FILTERTOOCLOSE(newTracks, activeTracks)
31:    activeTracks  $\leftarrow$  activeTracks  $\cup$  newTracks
32:    allTracks  $\leftarrow$  allTracks  $\cup$  newTracks
33:  end for
34:  return allTracks
35: end function
36:
37: function WINDOWBA(startPoses, images, tracks) ▷ See paper Section 3.2.4.
38:  currentWindow  $\leftarrow$  CREATESCENE(startPoses)
39:  windows  $\leftarrow$  {currentWindow}
40:  while last image not processed do
41:    currentWindow  $\leftarrow$  REMOVEUNCONFIDENTPOSES(currentWindow)
42:    baPoses  $\leftarrow$  SELECTBASUBSET(currentWindow)
43:    lastConfidentPose  $\leftarrow$  null
44:    while true do
45:      candidate  $\leftarrow$  PICKIMAGELININCREASINGOFFSET(images)
46:      if ISCONFIDENT(baPoses, candidate) then ▷ Includes optimizing for the candidate's pose against baPoses
47:        lastConfidentPose  $\leftarrow$  candidate
48:      else
49:        break
50:      end if
51:    end while

```

---

---

```

52:   currentWindow ← ADDPOSE(currentWindow, lastConfidentPose)
53:   windows ← windows ∪ {currentWindow}
54:   constraints = EXTRACTCONSTRAINTS(windows)
55:   constraints = FILTERBYPOSES(currentWindow.cameraPoses)
56:   currentWindow ← INITPOSESLINEAR(constraints)
57: end while
58: return windows
59: end function
60:
61: function ANCHORHANDLING(windows, costMatrix) ▷ See paper Section 3.3.
62:   anchorWindows ← ∅
63:   varianceMatrix ← ESTIMATEVARIANCES(windows)
64:   samplingMatrix ← (1 − costMatrix) ∘ varianceMatrix
65:   while |anchorWindows| < maxGlobalAnchors do
66:     cameraPair ← IMPORTANCESAMPLEDPAIR(samplingMatrix)
67:     window ← BESTWINDOWCONTAINING(windows, cameraPair.first)
68:     anchorWindow ← ADDPOSE(window, cameraPair.second)
69:     anchorWindow ← OPTIMIZE(anchorWindow)
70:     if ISSTABLE(anchorWindow, cameraPair.second) then
71:       anchorWindows ← anchorWindows ∪ {anchorWindow}
72:     end if
73:   end while
74:   return anchorWindows
75: end function
76:
77: function OPTIMIZE(scene)
78:   scene.points ← FINDPOINTS(scene.cameraPoses)
79:   scene.points ← SUBSAMPLEPOINTS(scene.points) ▷ See paper Section 3.2.2.
80:   scene ← BA(scene)
81:   while config says so do
82:     scene.points ← FINDPOINTS(scene.cameraPoses)
83:     scene.points ← FILTERBADPOINTS(scene.points)
84:     scene.points ← SUBSAMPLEPOINTS(scene.points) ▷ See paper Section 3.2.2.
85:     scene ← BA(scene)
86:   end while
87:   return scene
88: end function

```

---

## 4. 5k high resolution video

We reconstructed a scene from very high resolution video with 5120x2700 pixels and 901 frames. Timings and a visualization of the reconstruction can be found in Figure 11.

Note that the major bottleneck in this reconstruction was the subsampling of points whose acceleration is left for future work.

## 5. Large, cooperative reconstruction

We reconstructed an office room from video footage with 14254 frames that consists of several video sequences recorded with three different cameras. We used a Go-

Pro Hero3, a budgeted consumer Canon S100 photo camera and a Sony camcorder and recorded all video sequences in FullHD with frame rates between 25 and 60 Hz.

Example images, a visualization of the reconstructed scene and camera tracks as well as timing information can be found in Figure 12.

Refer to Figure 13 for an illustration of the effect of different framerates on the camera tracks.

### 5.1. Video

Please refer also to the supplied video for a more three-dimensional impression of the reconstruction.

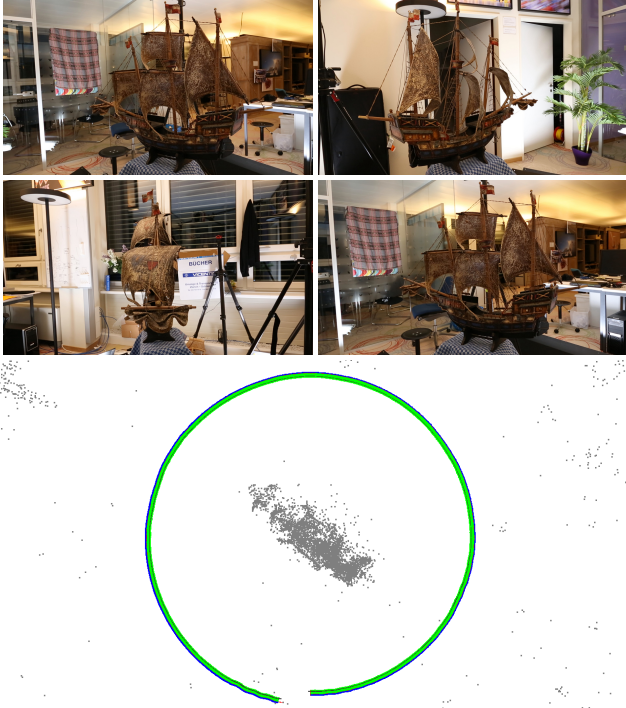


Figure 2. Ship scene: Input images and visualization - topdown view. Note that the loop is not closed because this scene was reconstructed without global anchor constraints to measure camera drift.



Figure 4. Bench scene: Input images and visualization - frontal view.

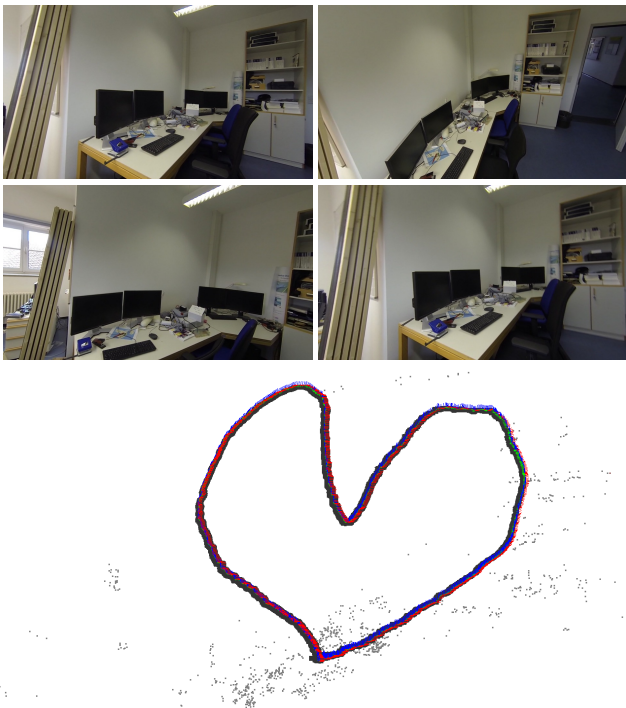


Figure 3. Office scene: Input images and visualization - viewing in camera direction.



Figure 5. Stairway scene: Input images and visualization - top-down view.



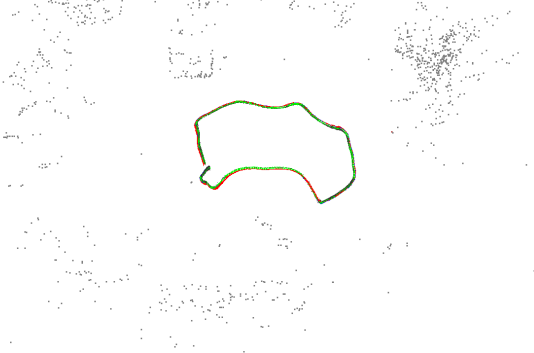


Figure 6. Office 2 scene: Input images and visualization - topdown view. Note that the loop is not closed because this scene was reconstructed without global anchor constraints to measure camera drift.

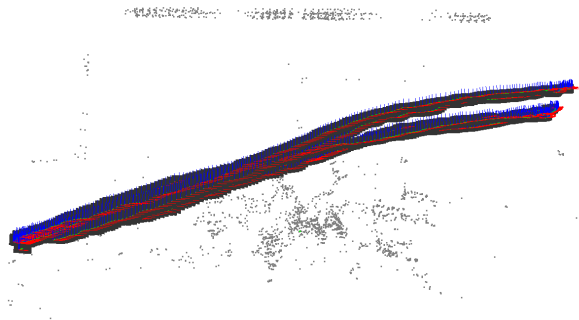


Figure 7. Swinging scene: Input images and visualization - viewing in camera direction.

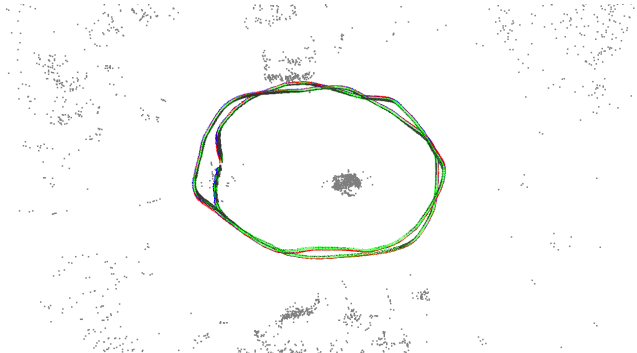


Figure 8. 2 loops scene: Input images and visualization - topdown view.

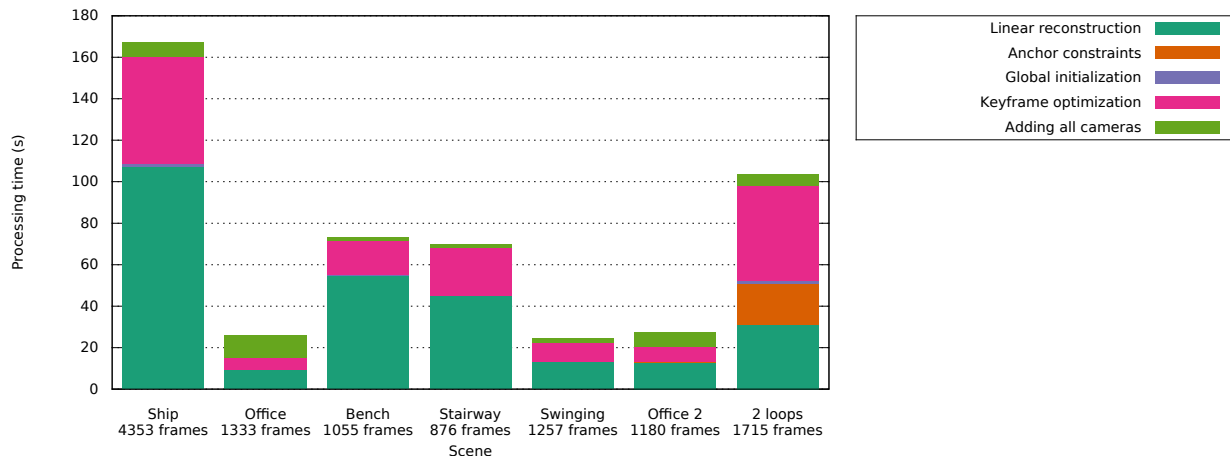


Figure 1. Reconstruction timings for the test scenes used in the paper. Reconstruction does not depend directly on the number of frames but on how much scene information was captured with a specific camera movement. Redundant or unnecessary information is skipped by our frame selection and point subsampling. Note that most scenes were reconstructed without anchor constraints to be able to evaluate the camera drift.

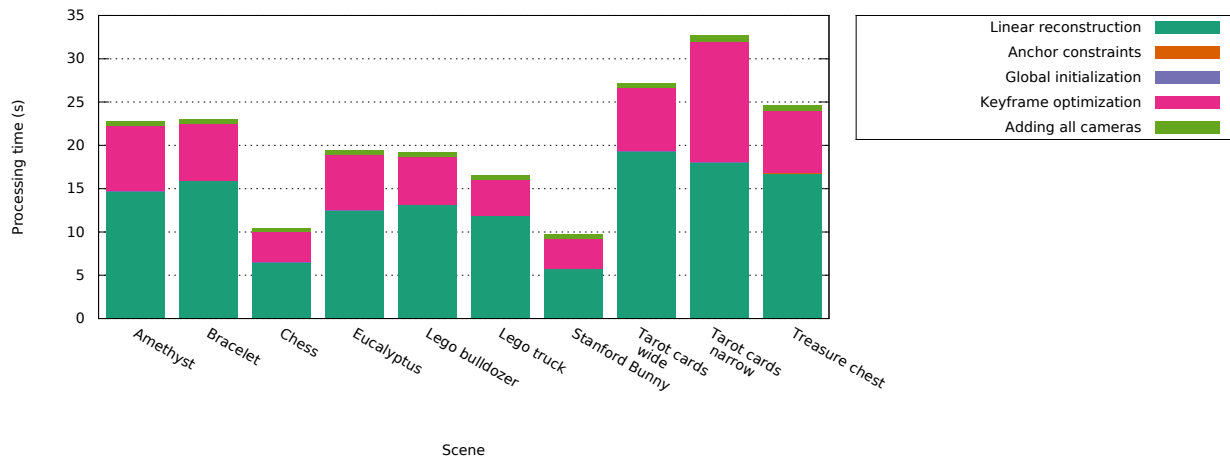


Figure 9. Timings for the stanford dataset reconstruction. Each scene contains 289 10MPixel images.

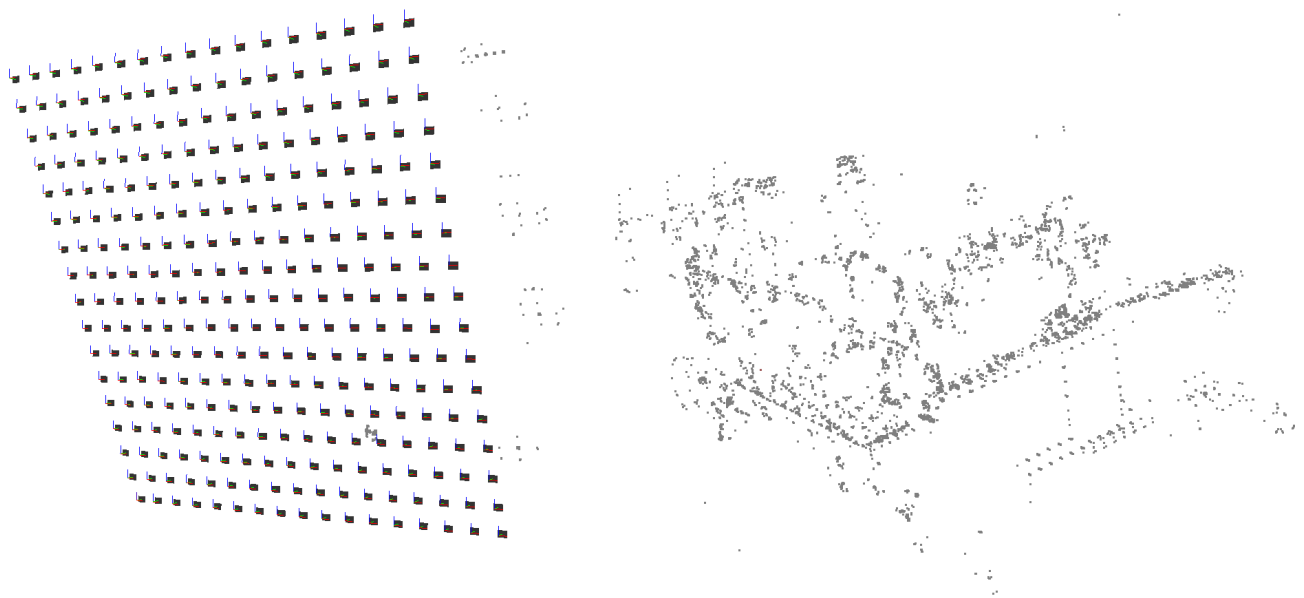


Figure 10. Reconstruction of the Stanford Light Field Archive lego bulldozer scene.

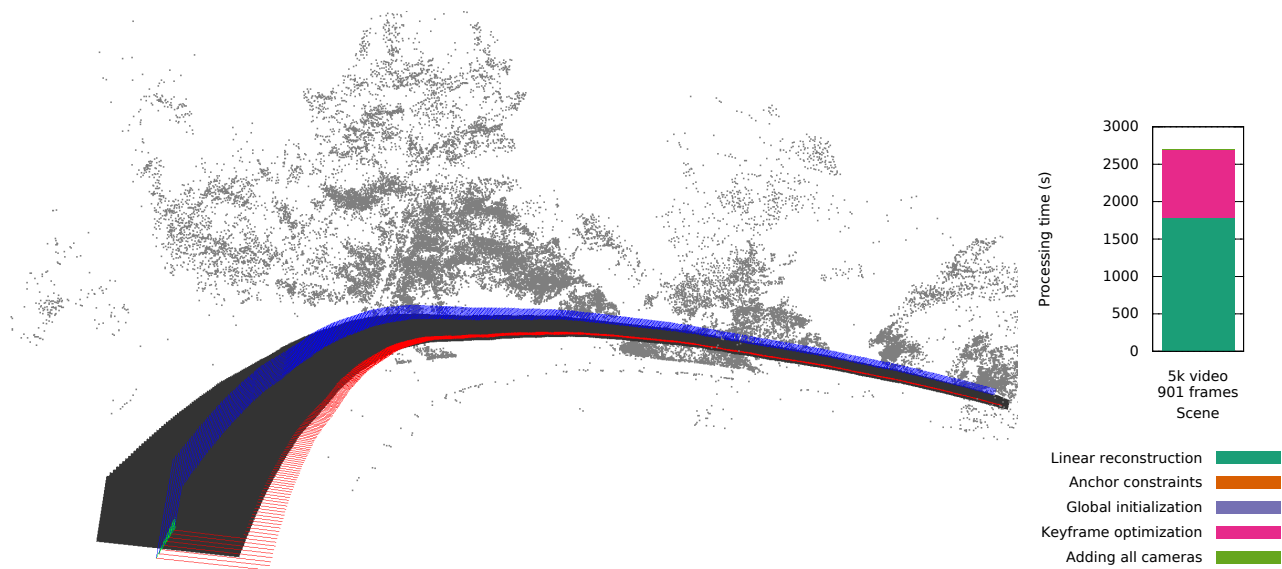


Figure 11. Reconstruction visualization and timings and reconstruction from a very high resolution (5k) video. Most time is spend on point subsampling during window BA and the final BA.



Figure 12. Example images, reconstruction visualization and timings from cooperatively captured video sequence set containing 14254 frames from three different cameras. Examples images from left to right: Canon S100, GoPro Hero3, Sony camcorder. All reconstructed camera trajectories are linked into a common global context by the anchor constraints.

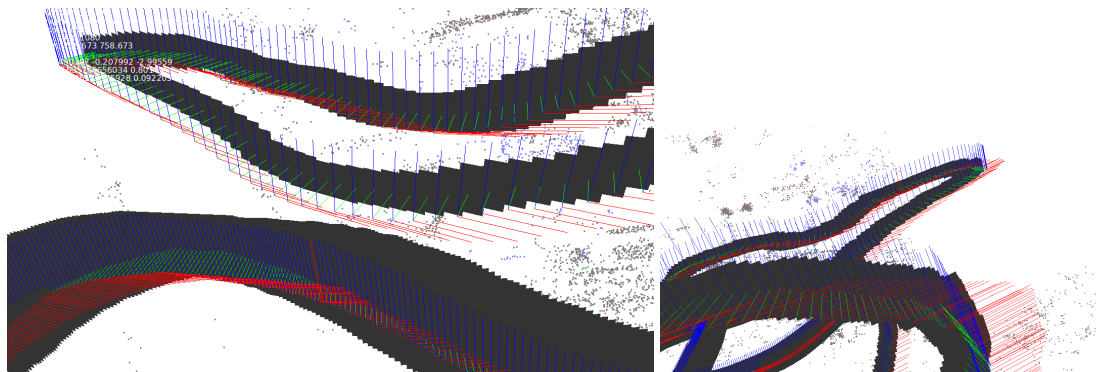


Figure 13. Effect of different framerates on camera trajectories. Left: Canon S100 (25 fps) vs. Gopro Hero 3 (60 fps). Right: Canon S100 (25 fps) vs. Sony camcorder (50 fps).