

Swipe Mosaics from Video

Appendix

Paper 1029

1. Baseline Algorithms

In this document we provide further results of comparisons between our system and Wu *et al.*'s VisualSfM [Wu07, WACS11], Viewfinder Alignment (VfA) of Adams *et al.* [AGP08], a “micro-SfM” method which computes a 2D translation using SIFT matching with RANSAC, and direct (all-pixel) methods as summarized by Szeliski [Sze06]. For each baseline algorithm we have the capability to load the output into our viewer for qualitative comparison. To load VisualSfM's output of 6d.o.f. camera positions into the Swipe Mosaic viewer we project each camera to a location and orientation on the 2D plane. We do this by computing a normal vector from the average “forward” direction over all cameras, then projecting each camera perpendicularly onto a plane defined by this normal. The location on this plane gives a 2D coordinate to be loaded into the viewer (see main paper). The inputs used are displayed Fig. 1.

1.1. Direct Methods

An excellent summary of techniques for Image Alignment and Stitching was presented by Szeliski [Sze06], Chapter 3, “Direct (pixel-based) alignment”, details a number of methods to compute a 2D alignment between image pairs. The general approach is to define some error metric which can evaluate how well each potential 2D alignment matches the contents of each image. Given the error metric, one can exhaustively evaluate all possible alignments, or use a coarse to fine method to limit the amount of computation, and the alignment which produces the lowest error is chosen. Various error metrics can be defined on overlapping pixels, such as Sum of Square Differences or Sum of Absolute Differences. We compare to the popular method of Normalized Cross Correlation (NCC), which we also used as the basis of our feature vector computation. NCC is an improvement over improves over Cross Correlation (which has a tendency to give incorrect offsets in the presence of large high intensity areas) by normalizing the overall intensity of each of the regions being compared. However, as noted by Szeliski, “[NCC's] performance degrades for noisy low-contrast regions” so the improved technique is not immune to problems.

Two temporally adjacent frames (Fig. 2a and Fig. 2b) from the PRISM sequence were selected. Some scene geometry is visible on the left of the image, but most of the pixels have been overloaded by the bright light and are reporting close to perfect white. Nevertheless, it is clear that horizontal camera motion has taken place,

given the parts of the geometry which we *can* see. The NCC image computed in MATLAB is shown in Fig. 2c. Two closeups of the region around the peak are shown in Fig. 2d and Fig. 2e. Note that the right hand peak in Fig. 2d has a higher NCC value, but the location (768, 432) implies that *zero* translation is the optimal image alignment. The left hand peak in Fig. 2e has a slightly lower magnitude, but the offset (679, 432) indicates a horizontal offset of $(768 - 679)/2 = 39$. Indeed, a translational shift of (39, 0) does bring the two images into good alignment. Our RRF based system produced an estimate with mean (0.0169, 0.0004) and variance (0.00132, 0.00168) for these images (note these results are not in units of pixels as above, so cannot be directly compared), showing primarily horizontal motion with roughly isotropic variance, as we would expect from the fact that the visible texture in the image confirms no vertical motion has taken place. Any algorithm which simply computes NCC over entire images, finds the single peak and uses the value will be prone to fail in image pairs such as this, whereas our system produces a translation in the correct direction. Note that un-normalized Cross Correlation for this image actually produced a *purely vertical* translation, reinforcing the idea that un-normalized cross correlation is unsuitable for large high intensity regions, performing even worse than NCC.

Our method returns a distribution over possible camera motion, which is a great advantage in cases of ambiguity such as repeated structure. Fig. 3 shows how in the presence of multiple potential alignments. Another example of how our method is superior to NCC is in shown in Fig. 3. When repeated structure creates a number of possible alignments, our method (Fig. 3d) returns an anisotropic estimate compared to the deterministic estimate provided by NCC (Fig. 3c). Corresponding repeated structure is shown in the NCC response image; the peak happens to be located in a ridge which indicates slight upward motion (as well as rightwards). The vertically adjacent ridges have similar NCC values and (given the magnitude of the camera motion between the images) are surely almost as likely. However the single transformation returned by pure NCC alignment will not represent this information at all. The result from our method is more desirable in this situation.

1.2. Structure from Motion: VisualSfM

VisualSfM [Wu07, WACS11] is a state of the art SfM system, which can process images either as “ordered” (temporally sequential frames) or “unordered”. It produces excellent results in general,

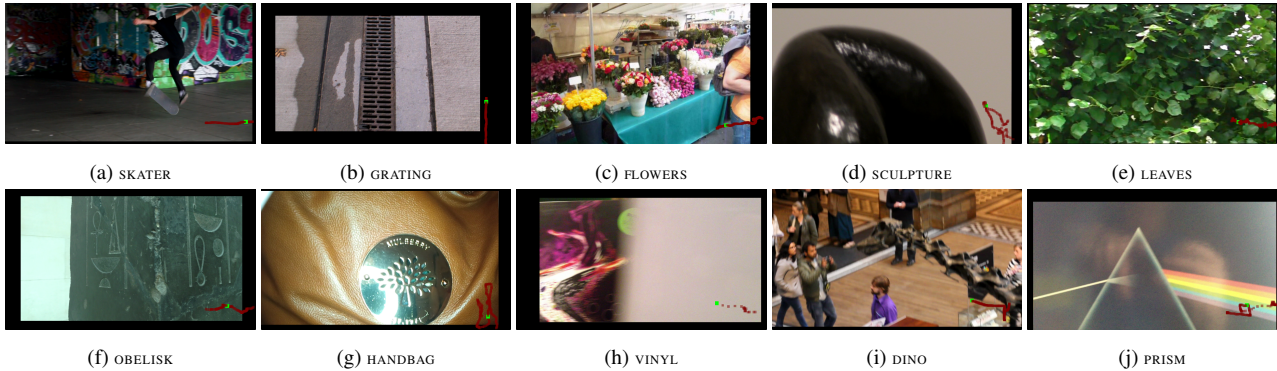


Figure 1: Screenshots of various sequences loaded in the Swipe Mosaic viewer

but struggles when few or misleading feature matches are present. We selected the *SCULPTURE* (Fig. 4a) and *LEAVES* (Fig. 4b) sequences as likely to cause SfM failure. *SCULPTURE* includes specularities, motion blur, and has few suitable corners for interest point detection. *LEAVES* contains lots of geometry suitable for interest point detection, but most of these areas are on leaves, which are being blown around in the wind, meaning points detected on them may adversely contribute to the optimization.

VisualSfM was run on each sequence in both ordered and unordered mode. This mode affects which image pairs are compared to find interest points; either all pairs (unordered) or only temporally adjacent pairs (ordered).

Running unordered on *SCULPTURE*, 38 camera locations were reconstructed from the 101 input images, leading to an incomplete Swipe Mosaic. Running ordered mode yields an even worse result, reconstructing 35 camera locations but in three independent groups. The locations that VisualSfM did produce were accurate, but the full camera path produced by our system is preferable, as shown in Fig. 5.

The *LEAVES* dataset consist of 201 images. VisualSfM in ordered mode computed locations for only 70 of the images, albeit producing a reasonable Swipe Mosaic. Running SfM in unordered mode computes a location for all of the images, but the placement undergoes a catastrophic failure, with the resulting Swipe Mosaic suffering severe artifacts. The failure takes the form of one side of the horizontal path being relatively correct, and the image locations gradually worsening as we travel along the video timeline, until the predicted image locations do not even overlap (see supplemental video). For both the sequences in this section, our system produced easily navigable locations for all cameras (see video and Fig. 4). We surmise that unordered mode producing better results in this case was due to a greater variety of image pairs being run through SIFT matching, rather than merely a few temporal neighbors. If a several consecutive frames of video are blurred or contain confusing motion, unordered mode will still be able to search for SIFT matches between frames “either side” of the problem area, thus providing a more robust solution.

1.3. μ SfM

“micro-SfM” or “ μ SfM” is a system which we have developed with the intention of it being an equivalent system to SfM, but without computing any structure, and with camera transforms limited to 2D translation. The thinking behind this is that computing a 6D quantity for each frame is an inherently harder task than computing a 2D quantity for each frame, and so simply comparing our 2D RRF method to VisualSfM was not a fair comparison. Rather, we should apply the technique from VisualSfM to the strictly easier problem of computing translations (not fundamental matrices) in order to compare like-for-like results. “ μ SfM” matches two images by generating SIFT descriptors and performing matching using the standard algorithm of Lowe [Low04]. A translation is computed using RANSAC to iteratively select a random SIFT match, compute the corresponding 2D transform and count the number of inliers. The transform with the highest inlier count is used to generate a final refined transform from the entire inlier set. To combine multiple translation estimates across an image sequence, we use the least-squares based layout algorithm developed for our RRF estimates. As the layout algorithm allows for a weighting to be applied to each relative transform (in our main system we use the inverse variance from the forest) we test two versions of μ SfM – one unweighted, and one weighted using the ratio of inliers to total number of matches found when generating the transform. This should ensure that translations for which every single SIFT match agrees are given more weight by the optimization.

Considering the simplicity of the method, μ SfM is surprisingly capable. In particular, it can produce just as good a camera path for the *LEAVES* sequence as our technique. However as the method relies entirely on interest point matching, we know it is susceptible to fail in the presence textureless regions, motion blur or repeated structure. The *VINYL* sequence contains a blurry obstruction which is very close to the camera, separating to regions containing strong texture information. The camera travels horizontally, starting in one textured region, passing the obstruction (which takes up the whole screen for a few frames) and ends viewing the second textured region. Surprisingly, inside the (apparently) textureless region, SIFT is able to detect a few interest points. Matching these interest points proves difficult however; most of them have extremely similar appearances, and despite implementing Lowe’s technique for avoid-

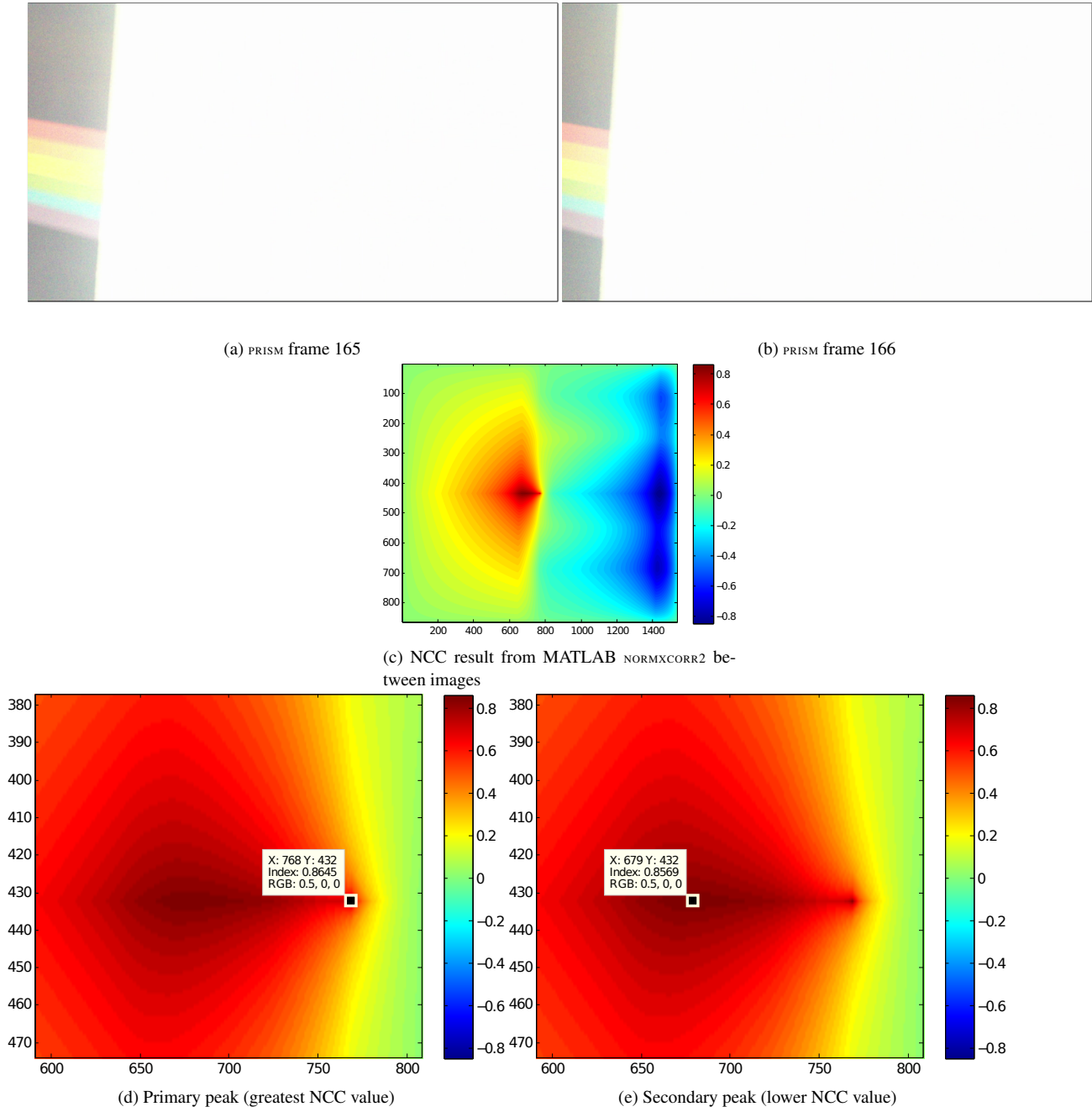


Figure 2: NCC failing to compute the correct offset for two frames in PRISM. Black borders added to top images for clarity.

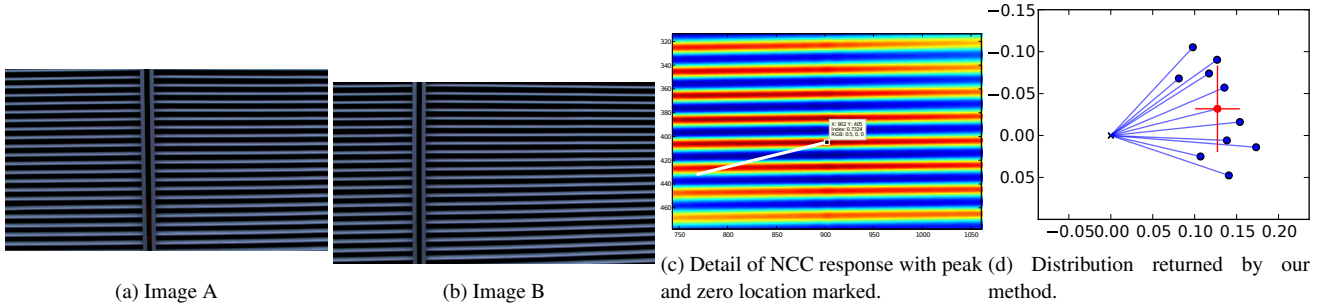


Figure 3: NCC computes a large diagonal offset for 2 frames containing repeated structure. The white line in c) connects the NCC peak with the location which would represent zero offset

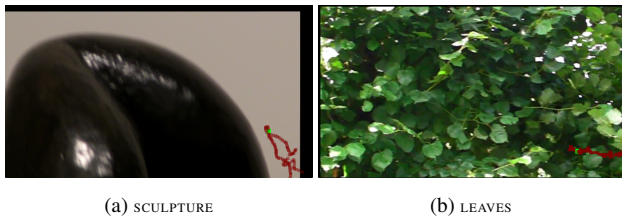


Figure 4: Screenshots in the Swipe Mosaic interface of the datasets on which we compare our performance to SfM. The minimap in the bottom right shows the camera locations.

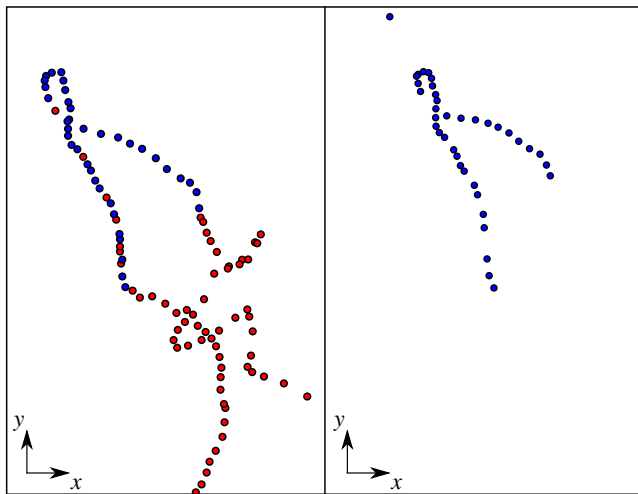


Figure 5: 2D camera coordinates (unitless) for SCULPTURE produced by our system (left) and VisualSfM (right). Blue points indicate images where both systems gave an estimate of location (note the estimates differ); red points are images where only our system produced an estimate. Note the obvious outlier at the top of the SfM result.

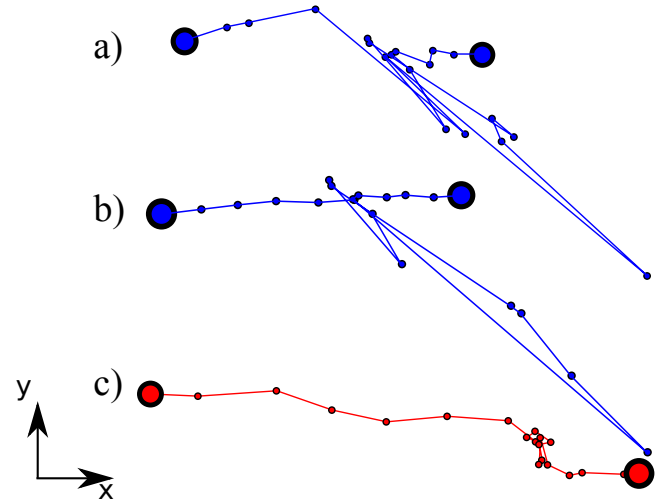


Figure 6: Regularized locations for the VINYL scene, for the images which straddle the obstruction. Dots represent image locations and the line joins them in temporal order. Larger dots show the start and end of this subsequence. As the correct camera motion is approximately a constant horizontal velocity, the ideal result would be equally spaced dots on a horizontal line. Scale between the 3 diagrams is not meaningful. a): μ SfM result with unweighted regularization. b): μ SfM with weighted regularization. c): our system.

ing ambiguous matches, the translations returned from the middle frames in this sequence were extremely noisy. Note the results in Fig. 6, bearing in mind the ideal answer would be almost pure horizontal motion. Both μ SfM results display problems with some frames ending up at a large displacement to the lower right corner of the map, with the subsequent frames on the normal timeline. It can be seen that the weighted version displays a smoother timeline at the beginning and end of the sequence, but for both a) and b) the mistakes in the middle of the sequence make this difficult to navigate in our interface (frames displaying the obstruction are incorrectly displayed amongst the frames of texture objects). By contrast our result, c), whilst by no means perfect, is a vast improvement on both μ SfM results. For the images where texture is available it computes a consistent horizontal motion. For the frames containing no texture, there is insufficient information to state which (if

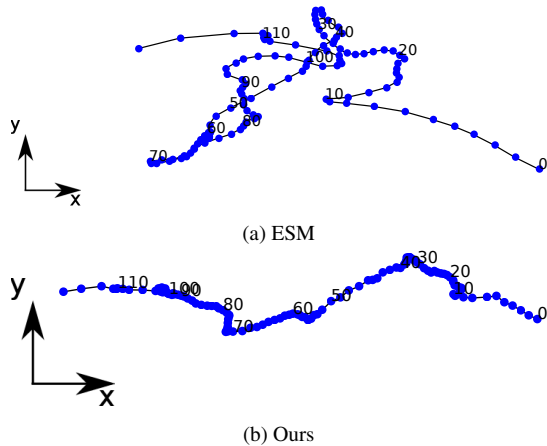


Figure 7: 2D camera coordinates produced for *OBELISK* using ESM (Fig. 7a) and our method (Fig. 7b). The true camera motion is approximately constant horizontal translation, coupled with rotation to remain pointed towards the object of interest (see supplemental video).

any) direction the camera has moved, so our system returns a number of zero mean, wide variance offsets. The optimization places these roughly on top of each other, generating the point cluster in the result. Obviously this is not actually correct, as the camera was always moving, but as there is no way to tell this simply from the images pairs we produce a reasonable result, which allows the sequence to be browsed as a Swipe Mosaic without artifacts. It may be possible to improve this aspect of the system by using a camera motion model in the layout algorithm, meaning that when we knew the first few frames had the camera move to the right, then when presented with insufficient visual information our estimate would be some kind of rightwards motion, rather than zero mean motion. Another failure case for μ SfM was *PRISM*, where the camera auto-gain causes whiteout for a few frames. Similarly to the previous sequence, no reliable feature matches could be detected during this central part of the image sequence, leading to incorrect matches in the middle of the sequence again.

1.4. Efficient Second-order Minimization

Efficient Second-order Minimization (ESM) [BM04] is a direct method which explicitly models the scene as a plane, and searches for a parameterized transform which minimizes the sum of squared differences between two images. The transform can be parameterized as anything from a full homography (8DoF) to a translational transform (2DoF), and the parameters are solved for using an efficient method which achieves Newton method like convergence rates, without having to compute the Hessian. We used the implementation of ESM available in Ed Rosten’s LibCVD project, using 2 DoF to produce a translation between each image pair, before using our layout algorithm. ESM produces a good Swipe Mosaic result for some of our test sequences, but the explicit parameterization of the scene as a plane leads to problems when faced with non-planar camera paths or distorting objects. ESM performed very badly on the *OBELISK* scene, laying out frames which should be very far from each other in roughly the same place. Our system produced

an intuitively navigable Swipe Mosaic. The camera paths produced by ESM and our method are shown in Fig. 7. Note the broadly horizontal linear path produced by our method in contrast to the ESM path which continually crosses itself. This is due to our method’s more gradual degradation as scenes deviate from planar, allowing us to cope with strong perspective deformations. See the supplemental video for a comparison of the browsing experience of these two solutions.

1.5. Viewfinder Alignment

Our final baseline comparison is to Viewfinder Alignment (VfA) of Adams *et al.* [AGP08]. VfA is a method to compute constrained transforms between temporally close video frames. VfA computes a “digest” for each frame by encoding edge information at 4 equally spaced orientations using gradient integral projection arrays, as well as detecting the top k peaks in the image. Two image digests are aligned by first calculating a single 2D shift which best aligns the edge information stored for each image. This 2D shift is applied to the detected corners of one of the digests. The number of inliers (a pair of points, one from each image, landing within 3 pixels from each other) between the two points sets is counted and taken as the confidence that the images have been aligned correctly. The set of inliers is used to generate a similarity transform, giving 4 degrees of freedom (translation, rotation and scale) between pairs of frames. Note that our VfA test scenes were chosen intentionally so that the rotation and scale change was negligible so these parameters are ignored, *i.e.* we are only interested in the relative accuracy of the *translations* computed by different methods. Experiments showed that when compute the scale change was typically between 0.98 and 1.02, and the rotation on the order of 0.1 radians, justifying this decision.

VfA has a number of attractive properties, including computational efficiency and being extremely resistant to noise. A disadvantage of the algorithm is that it is completely deterministic, in that only one 2D translation between each frame pair is considered, when the digest edge information could be used to produce a distribution of translations. Additionally, the corners returned from the corner detector are simply stored as 2D locations, without any kind of descriptor, allowing corners which represent different scenes points to potentially be aligned with each other and treated as an inlier. We compared our system to our own re-implementation of VfA. This code is supplied as supplemental material. We now present detailed analysis of VfA on various test scenes.

All the sequences in this document were run through our re-implementation of Viewfinder Alignment. We tried to match each digest with the digests from other images which were within 6 frames (forward or backwards). If a complete graph could be constructed, we used all the inferred translation values as input to our linear least squares regularization (see main paper). Inlier matrices are shown using the standard Jet colormap, except that pairs which either produced zero inliers or were not compared (*i.e.* they were too far apart temporally) are left blank.



Figure 8: Sample Images from lobby sequence.



Figure 10: Sample Images from fence sequence.

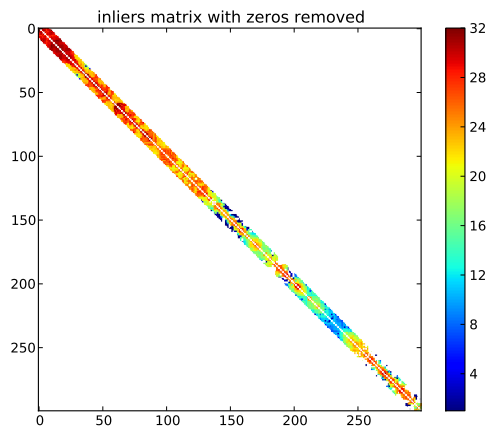


Figure 9: Number of inliers from lobby sequence.

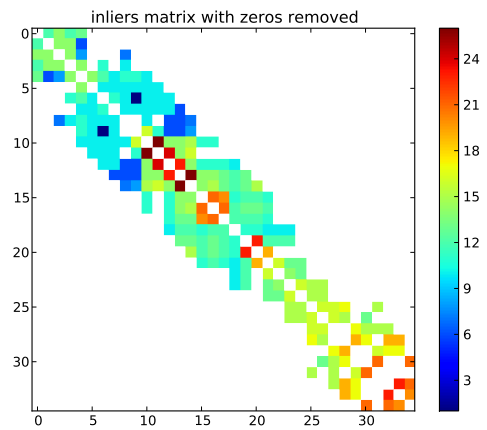


Figure 11: Number of inliers from fence sequence.

Successes

Lobby

The lobby sequence is largely featureless, but the objects that are seen display strong vertical and horizontal edges, resulting in an excellent output when browsed as a Swipe Mosaic.

Fence

The fence contains repeated structure with many similar looking horizontal edges but VfA is robust to this, producing a fully connected set of inliers (Fig. 11) and a correct reconstruction.

Failure cases

Vinyl

The Vinyl sequence contains an obstruction which does not trigger Viewfinder Alignment’s corner detection (Fig. 12). This causes a large region with zero inliers (Fig. 13) which prevents the start and end of this sequence from connecting to each other, and therefore renders it impossible to create a single set of camera paths to be loaded into our viewer.

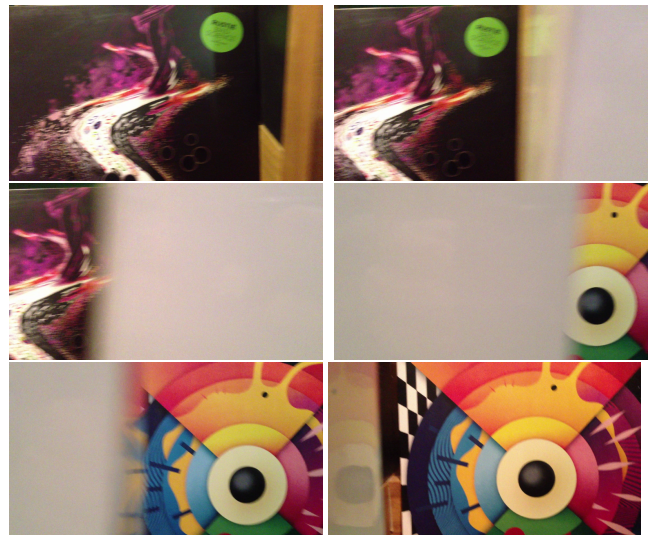


Figure 12: Sample images from vinyl sequence.

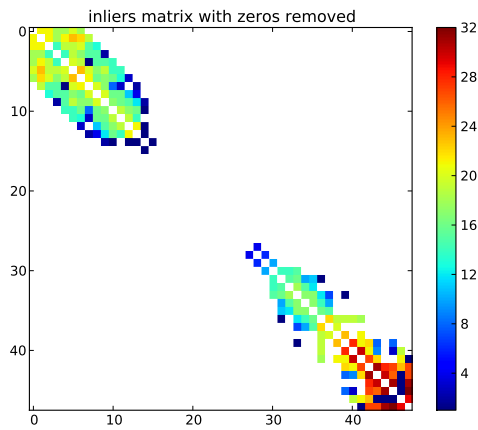


Figure 13: Number of inliers from vinyl sequence. Note two separate “islands”.

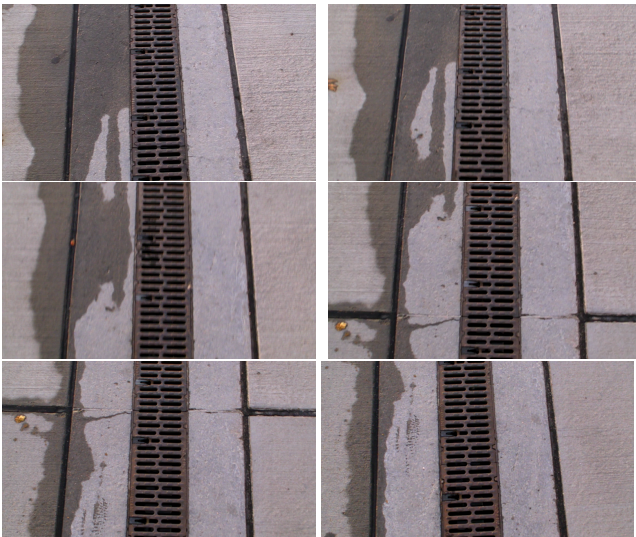


Figure 14: Sample images from `grating` sequence.

Grating

The Grating sequence is an interesting case because it contains very easily localisable vertical edges but relatively few unambiguous horizontal edges (Fig. 14). Viewfinder Alignment finds sufficient transforms to regularise the 80 frame segment all together as one connected cluster, leading a promising looking inliers graph (Fig. 15). However, a small number of incorrect matches corrupt the whole regularisation, resulting in final camera locations shown in Fig. 16 and Fig. 17. The correct arrangement should be a roughly straight vertical line. It seems likely that the lack of strong horizontal edges meant that when an incorrect alignment was proposed and approved by the corner correspondence stage of VfA. Corners are deemed as inliers based on whether a given shift puts them on top of each other, not based on any kind of descriptor based on the visual appearance of whatever was originally detected as a corner.

submitted to *Vision, Modeling, and Visualization* (2016)

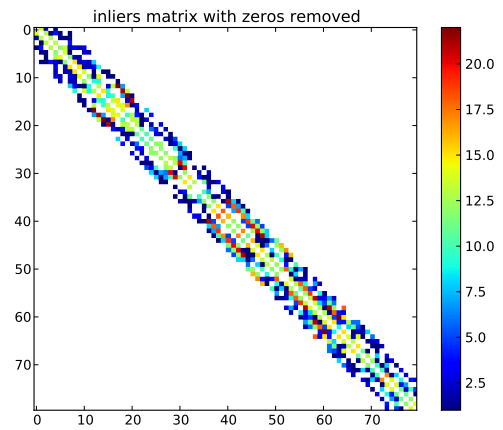


Figure 15: Number of inliers from grating sequence.

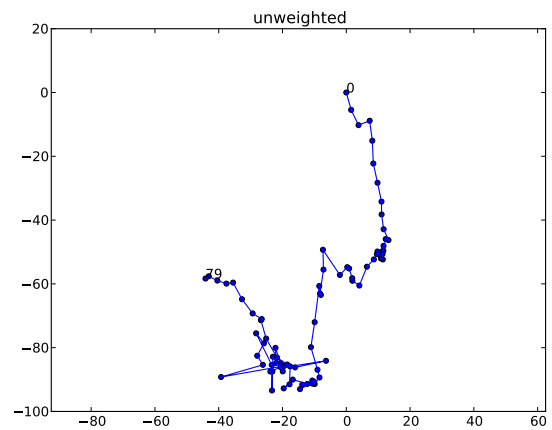


Figure 16: Regularized camera locations for grating, unweighted. 0 is the first frame; 79 the final frame.

In this situation potentially even adding corner descriptors to the algorithm would not remedy the situation, because it seems like most corners are liable to be detected on either the drain or the two grooves next to it, and any hypothetical descriptor computed on these locations is likely to be visually similar to another descriptor computed somewhere else on the drain / groove. The best movement cues in this scene are the water stains on the floor, which are non-repeating, but these are not captured well by the VfA digest.

2. Further Swipe Mosaic Results

2.1. Synthetic Satellite footage dataset

As well as the real video sequences, we converted a timelapse video of the Earth recorded from the International Space Station into a Swipe Mosaic. We first constructed an intermediate video by cropping out a thin horizontal strip from the bottom of each image and splitting each strip into eight overlapping images. A virtual cam-

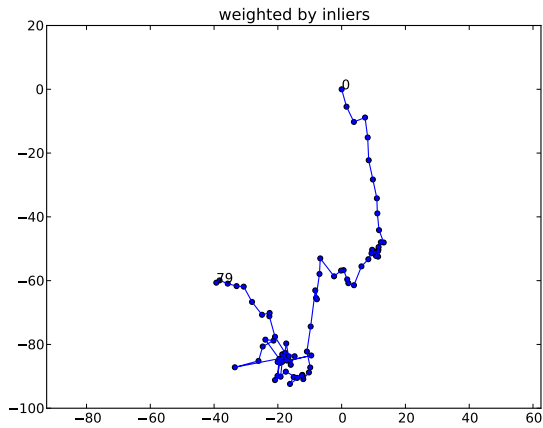


Figure 17: Regularized camera locations for grating, weighted by inlier count. 0 is the first frame; 79 the final frame.

era was moved back and forth along the strips, moving forward in time upon reaching the end, to give these cut up frames a nominal temporal ordering. We run pairwise prediction, regularization and translational loop closure on this sequence, and know that the images should ideally be estimated to form a regular rectangular grid.

It is hoped each image would know from the translational RRF that their neighbors on the same strip were at a purely horizontal offset, and neighbors on a different strip were at a purely vertical offset. The output of the first regularization step is shown in Fig. 19a. The arrangement is approximately what we would have hoped for, but the locations as a whole “lean” to one side. This can be explained by noting that for the images at either end of the strip, when the virtual camera moves “up” or “down”, the overlapping pixel data between the image at either end of this link will actually move diagonally, because all the earth’s surface appears to be moving away from the focus of expansion. Because the strips were not symmetrically cropped (the main goal of the cropping was to remove the visible parts of the ISS which appeared in the frame, of which there was more on one side) we see that the “up-down” links such as (7, 15) push the entire system to the left to a greater extent than the links on the other side such as 40, 48.

Ideally, our loop closure step should (with slightly modified thresholds to account for the shorter loops present in this artificial scene versus a real scene) detect loop points between every image and its corresponding vertical neighbors, *i.e.* link (0, 15), (1, 14), (2, 13) *etc.*. After incorporating a pairwise prediction from each and re-regularizing, we would hope to see the same overall grid structure, but with less of the horizontal skew visible in Fig. 19a. The result of automatic loop closure is shown in Fig. 19b. Loop points have been found in the majority of places we hoped to see them. The output is not perfect but the transitions between, for example, images 3, 12 and 19 are closer to vertical than before the loop closure, and so are correspondingly improved in the visualization. For this dataset, vertical would be the ideal answer. Even

with the imperfect results, the iss sequence is easy to navigate as a Swipe Mosaic.

3. Implementation Details

3.1. Training Data Acquisition

Capturing real-world video data with ground-truth camera motion is error-prone and time consuming even with specialized equipment. After a variety of attempts, including using multi-camera rigs and improvised outdoor motion capture, we eventually chose to generate synthetic image pairs with known camera motion. The RRFs are able to learn how different 2D translations and rotations appear when the world is shiny, smooth, bumpy, repetitive, and when distracting objects are moving about. We did not render motion blur, but this is certainly possible. Synthesizing training data with graphics techniques has previously proved successful [SFC*11, MAHPB12], despite the obvious risk that the resulting regressor or classifier may only be accurate on artificial-looking scenes. Aiming for large variations in shape and appearance, we rendered a family of random landscapes consisting of both angular pillars and smooth NURBS surfaces, with shape variability generated by randomly moving the pillars and deforming the surfaces. Appearance variability was achieved by rendering each object with a random color and reflectivity. We render two images of each landscape, with a random in-plane camera translation as the only difference between them. The generated images include both texture rich and texture poor regions, and irregular curved edges between NURBS surfaces, which are elusive to many interest-point detectors.

A benefit of our supervised learning approach is that if deficiencies are found in the future, it is possible to augment the training set and improve model performance. During development of our system, it was determined that regularly repeating structures posed difficulties for the system. We augmented the training set by adding “billboard” datasets which replaced the random landscape previously described with a textured polygon, containing one of a set of images of repeated structure which were obtained from Flickr and other Creative Commons sources. Example frame pairs from our training data are shown in Fig. 20.

3.2. Gabor Filter Bank

As mentioned in the the *Feature Computation* section of the main paper, a bank of Gabor filters are used as part of the feature computation process. Each filter is computed from the product of a Gaussian and a sinusoid, according to (1).

$$g(x, y; \lambda, \theta, \sigma, \gamma) = \exp\left(-\frac{\hat{x}^2}{2\sigma^2} - \frac{\hat{y}^2}{2\sigma_y^2}\right) \exp\left(\frac{2\pi\hat{x}}{\lambda}\right) \quad (1)$$

$$\hat{x} = x \cos \theta + y \sin \theta \quad (2)$$

$$\hat{y} = -x \sin \theta + y \cos \theta \quad (3)$$

$$\sigma_y = \frac{\sigma}{\gamma} \quad (4)$$

$$(5)$$

λ represents the wavelength of the sinusoid, θ is the orientation of the sinusoid (the orientation parameters allows the detection of

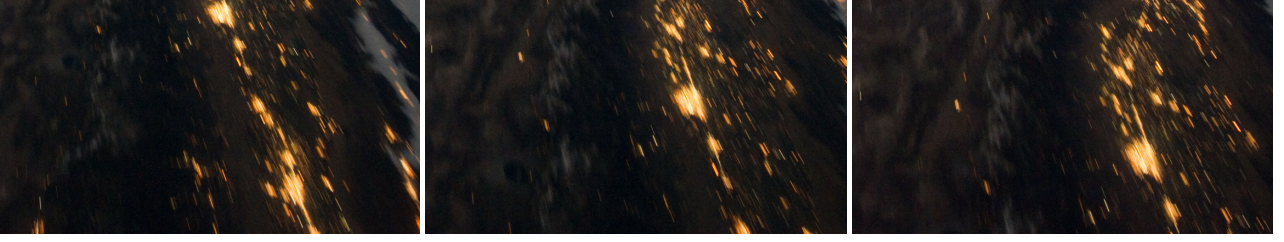
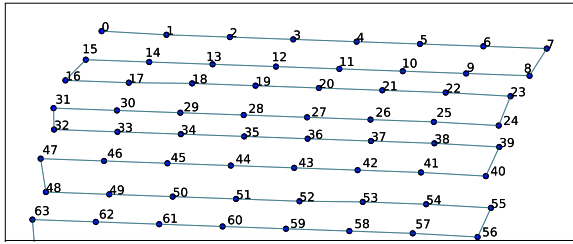
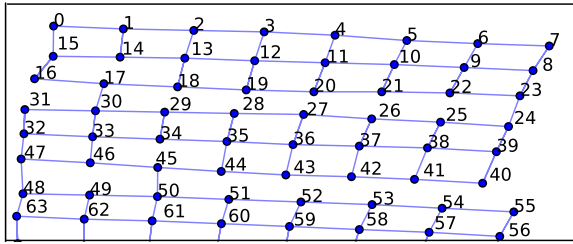


Figure 18: Example images from ISS sequence, from the right hand edge of the image grid. Note how structures travel diagonally down to the right between frames – this is due to the curvature of the earth from the ISS’ vantage point.



(a) Without loop closure.



(b) With loop closure.

Figure 19: Loop closure on iss. For this sequence only, the loop point detection parameter was set to 10 to allow for shorter loops. Only locations of the initial 64 frames are shown for clarity.

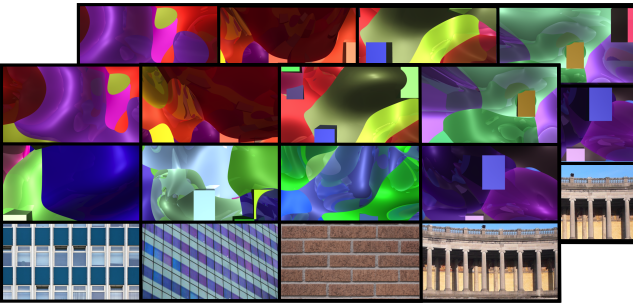


Figure 20: The RRF is trained on thousands of two-frame image sequences, with known camera transforms. To obtain sufficient quantity and variety of camera moves and scenes, we generated the training data using a custom-built but simple graphics engine. The top 2 rows show a few examples of the procedurally generated scenes with depth variations and dynamic scene elements. The bottom row contains real images from Flickr that were mapped onto flat moving surfaces to generate training data with repeated textures.

λ	θ	σ	γ
100	0	4	1
10	$0, \frac{\pi}{4}, \dots, \frac{7}{4}\pi$	2	1
10	$0, \frac{\pi}{4}, \dots, \frac{7}{4}\pi$	2	0.5
10	$0, \frac{\pi}{4}, \dots, \frac{7}{4}\pi$	3	1
10	$0, \frac{\pi}{4}, \dots, \frac{7}{4}\pi$	3	0.5

Table 1: Parameters used to generate the Gabor filter bank.

multimodal ridges at different angles), σ represents the standard deviation of the Gaussian, and γ controls how this standard deviation varies in the x and y directions (*i.e.* creating an elliptical function). The ranges of values used for these parameters is specified in 1. For each configuration of parameters, the filter is created as wide (in pixels) as necessary to encompass 3 standard deviations for the Gaussian.

3.3. Pseudo-code

```

procedure FEATEXTRACTIMG( $I_s, I_t$ )
  for  $l \in \{1, 2, 4, 6, 8\}$  do
    for  $x \in \{0 \dots l - 1\}$  do
      for  $y \in \{0 \dots l - 1\}$  do
         $x, y \leftarrow \text{TEMPLPIX}(l, x, y, I_t.\text{shape})$ 
         $\hat{x}, \hat{y} \leftarrow \text{SEARCHPIX}(l, x, y, I_s.\text{shape})$ 
         $T \leftarrow I_t[y, x]$ 
         $S \leftarrow I_s[\hat{y}, \hat{x}]$ 
         $N \leftarrow \text{NCC}(S, T)$ 
        ENCODENCC( $N, l, S.\text{shape}$ )

```

```

procedure ENCODENCC( $N, l, \text{shape}$ )
  OUTPUT(MIN( $N$ ), MAX( $N$ ), MEAN( $N$ ))
   $x, y \leftarrow \text{PEAKCOORDS}(N)$ 
   $x', y' \leftarrow \text{NORMPEAK}(x, y, \text{shape})$ 
  OUTPUT( $x', y'$ )
  for  $p \in \{10, 20\}$  do
    OUTPUT(LAPLACECOORDS( $N, x, y, p$ ))
  OUTPUT(NORMEDHIST( $N, (-1, 1), 5$ ))
  if  $l \leq 2$  then
    for  $G \in \mathcal{G}$  do
       $H \leftarrow N * G$ 
       $a, b \leftarrow \text{MIN}(H), \text{MAX}(H)$ 
       $c, d \leftarrow \text{MEAN}(H), \text{MEDIAN}(H)$ 

```

OUTPUT(a, b, c, d)

Pseudocode to extract a feature vector from equal sized grayscale images I_s and I_t . Indentation denotes structure, as with Python. Zero based indexing is used. $I[y, x]$ is a slicing operation to extract the subwindow defined by pixel index vectors x and y . TEMPLPIX returns the pixel indices to extract a small “template” window for a given level, window index and image size. SEARCHPIX operates similarly but produces indices for a larger “search” window. OUTPUT appends a variable number of features to the feature vector being built for this image pair (for clarity, there is no variable to represent the feature vector in the pseudocode). PEAKCOORDS computes the 2D pixel shift for the peak of the NCC response. NORMPEAK normalises this in relation to the template image size. LAPLACECOORDS computes a Laplacian coordinate descriptor on N around the peak point (x, y) at a scale p . NORMEDHIST returns a 5 bin histogram of the NCC image with bin limits $(-1, 1)$.

3.4. Performance

All results were generated by a single trained translational RRF with 10 trees and maximum depth 12. At each node, 2000 possible feature splits were considered. The one-time training took 1h 43m on a Core 2 Duo 2.8 GHz, using both CPU cores. All experiments were carried out with the same synthetic training set built from 8800 image pairs. The rotational RRF was trained with 20 trees of depth 12, considering 100 feature splits per node. Only 400 image pairs were used in the training set, as inferring pure rotation is a strictly easier problem than translation, because what appears in the images is not scene geometry dependent. All images were 768×432 or 640×480 . A 2.5 GHz single core Xeon, computed 121 NCC matches and the 3599D features for two images in 15 seconds. Most useful datasets contain thousands of image pairs, so we used a cluster to process datasets quickly.

3.5. Optical Flow

Initial versions of our feature extraction used Optical Flow (rather than NCC) on the input images, before condensing that information into a vector. Building the feature vector from NCC is not an obviously better choice than using Optical Flow, but we obtained better test-time results with NCC based features. It is likely that the NCC responses are better correlated with motion-confidence than flow, which has some estimated vector for every pixel. Significantly, the regressor has the benefit of learning from our graphics engine: it has seen thousands of rendered examples of image pairs, with knowledge of the true 2D Euclidean transform.

4. 6D to 2D Camera Projection

To compare the 6D ground truth poses with each 2D solution, three steps had to be carried out. First, a ground truth pose must be established for each RGB frame, as the FREIBURG2 dataset only provides camera locations from motion capture, and the Kinect and motion capture systems were running unsynchronised at different frequencies. Secondly, some projection of each 6D camera pose onto 2D must be established, and finally the 2D solutions must be scaled and aligned to assess the correctness of the locations.

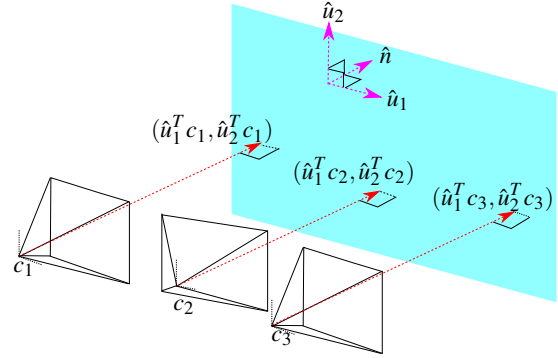


Figure 21: Diagram showing how 3D camera positions c_1, c_2, c_3 are projected (red dashed line) onto a 2D coordinate space defined by a plane (blue). For the “best fit” case, the plane normal \hat{n} is generated by averaging the forward vectors of all the cameras. The basis vectors \hat{u}_1 and \hat{u}_2 lie within the plane and are mutually orthogonal. When fitting to the plane defined by an individual camera, \hat{u}_1 and \hat{u}_2 are chosen to be parallel to the “right” and “up” vectors of individual cameras (black dotted lines).

To establish a 6D ground truth location for each RGB frame, we linearly interpolated between the closest two motion capture positions using the globally synched timestamps (available for both motion capture and Kinect readings). The rotation was encoded as a quaternion during this process, ensuring linear interpolation is a reasonable approach. The result of this operation is a 6D camera pose corresponding exactly to each RGB frame. The 6D to 2D projection step is described below, and the final alignment step is carried out using the Procrustes algorithm [Pri12].

Given a 3D plane represented as a unit normal \hat{n} and a point on the plane p , we define two unit vectors \hat{u}_1 and \hat{u}_2 such that $\hat{u}_1^T \hat{u}_2 = \hat{u}_1^T \hat{n} = \hat{u}_2^T \hat{n} = 0$. The exact orientation of these vectors is not important as the subsequent alignment includes a rotation step. Each camera location c_i is projected onto a 2D location on the plane $(\hat{u}_1^T c_i, \hat{u}_2^T c_i)$ (see Fig. 21). To define the plane orientation, we tried both a “best fit” to all the camera locations, and also tried using the “up” and “right” vectors of each individual camera in turn. The “best fit” plane was defined by setting \hat{n} to the average forward direction of each camera. The maximum angular difference between the computed normal and the forward vector of any of the cameras is 9.46 degrees, confirming that the scene contains only minimal rotation and is therefore a good candidate sequence for this comparison.

Additionally, Fig. 22 contains the comparison between our method and NCC for the FREIBURG2 dataset.

References

- [AGP08] ADAMS A., GELFAND N., PULLI K.: Viewfinder alignment. *Eurographics* (2008). 1, 5
- [BM04] BENHIMANE S., MALIS E.: Real-time image-based tracking of planes using efficient second-order minimization. In *Intelligent Robots and Systems* (2004), vol. 1, pp. 943–948. 5
- [Low04] LOWE D. G.: Distinctive image features from scale-invariant keypoints. *IJCV* (2004). 2
- [MAHPB12] MAC AODHA O., HUMAYUN A., POLLEFEYS M., BROS-

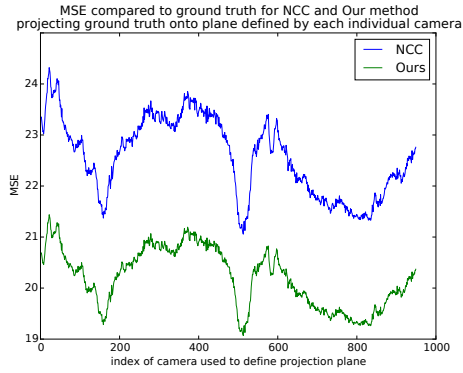


Figure 22: Comparison of Mean Squared Errors between our method and NCC for the whole range of possible individual camera projection planes, for FREIBURG2. Note the self similarity of the lines. As we change which ground truth camera defines the projection plane, both solutions (being so similar) have coinciding increases or decreases in performance.

- TOW G. J.: Learning a confidence measure for optical flow. *PAMI* (2012). 8
- [Pri12] PRINCE S.: *Computer Vision: Models Learning and Inference*. Cambridge University Press, 2012. 10
- [SFC*11] SHOTTON J., FITZGIBBON A., COOK M., SHARP T., FINOCCHIO M., MOORE R., KIPMAN A., BLAKE A.: Real-Time Human Pose Recognition in Parts from a Single Depth Image. *CVPR* (2011). 8
- [Sze06] SZELISKI R.: Image alignment and stitching: a tutorial. *Found. Trends. Comput. Graph. Vis.* (2006). 1
- [WACS11] WU C., AGARWAL S., CURLESS B., SEITZ S.: Multicore bundle adjustment. *CVPR* (2011). 1
- [Wu07] WU C.: SiftGPU: A GPU implementation of scale invariant feature transform (SIFT). cs.unc.edu/~ccwu/siftgpu, 2007. 1