

Anli Lim · Bharath Ramesh · Yue Yang · Cheng Xiang · Zhi Gao ·  
Feng Lin

# Real-Time Optical flow-based Video Stabilization for Unmanned Aerial Vehicles

Received: date / Revised: date

**Abstract** This paper describes the development of a novel algorithm to tackle the problem of real-time video stabilization for unmanned aerial vehicles (UAVs). There are two main components in the algorithm: (1) By designing a suitable model for the global motion of UAV, the proposed algorithm avoids the necessity of estimating the most general motion model, projective transformation, and considers simpler motion models, such as rigid transformation and similarity transformation. (2) To achieve a high processing speed, optical-flow based tracking is employed in lieu of conventional tracking and matching methods used by state-of-the-art algorithms. These two new ideas resulted in a real-time stabilization algorithm, developed over two phases. Stage I considers processing the whole sequence of frames in the video while achieving an average processing speed of 50fps on several publicly available benchmark videos. Next, Stage II undertakes the task of real-time video stabilization using a multi-threading implementation of the algorithm designed in Stage I.

## 1 Introduction

In recent times, unmanned aerial vehicles (UAVs) are often equipped with streaming video cameras that can be employed for immediate observation. They are popular for several applications such as rescue, surveillance, mapping, etc. The main drawbacks of videos taken by UAV

Anli Lim, Bharath Ramesh, Yue Yang and Cheng Xiang  
Department of Electrical and Computer Engineering  
National University of Singapore  
Singapore 117576  
Tel.: +65 65164258  
Fax: +65 67791103  
E-mail: bharath.ramesh03@u.nus.edu

Zhi Gao and Feng Lin  
Temasek Laboratories  
National University of Singapore  
Singapore 117576

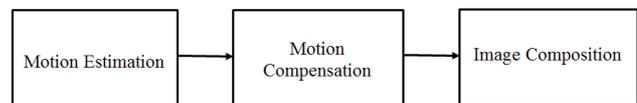


Fig. 1: Video Stabilization Framework

are the undesired shaking motion caused by atmospheric turbulence and jittery flight control of the platform. The shaky motion in the video proves to mitigate the fundamental aim of using UAV for vision-based tasks. In particular, the unstable motion also inhibits higher level vision tasks, such as detecting and tracking of objects in the video. While significant work has been done for handheld cameras and ground vehicles, there are limited works for stabilization of videos taken from UAV in the literature. Moreover, the existing UAV video stabilization algorithms, cannot be applied on-board UAV systems in real-time.

In general, UAV video stabilization algorithms (Shen et al (2009); Vazquez and Chang (2009); Wang et al (2011); Dong et al (2014); Lucas et al (1981)) follow three main steps: (1) motion estimation, (2) motion compensation and (3) image composition as indicated in Figure 1. Most methods revolve around finding the 2D motion model (such as homography) to estimate the global motion trajectory. Then a low-pass filter is applied to the trajectory to sieve out the high frequency jitters. Then the low frequency parameters are then applied onto frames via warping. This framework is really effective for scenes with very little dynamic movement which is applicable to aerial videos taken by UAVs. However, estimating the global motion trajectory using a window of frames cannot achieve instantaneous stabilization, as required for real-time processing.

In Shen et al (2009), a video stabilization algorithm for UAV was proposed using a circular block to search and match key places. The estimated affine transform was then smoothed by the polynomial fitting and pre-

diction method (PFPM). However, this approach only achieved a speed of less than 10fps on a desktop with 3.0GHz processor and 1GB RAM for images with resolution of 216x300p.

In Vazquez and Chang (2009), a smoothing method utilizes Lucas-Kanade tracker(Lucas et al (1981)) to detect interest points. The unintended motion compensation was accomplished by adjusting for extra rotation and displacements that generate vibrations. This approach is able to achieve a stabilizing speed between 20fps and 28fps for images with resolution of 320x240 pixels on a MAC laptop with 2.16GHz Intel Core 2 Duo Processor with a three frame delay.

Wang et al (2011), later proposed a three-step video stabilization method for UAVs. Firstly, a FAST corner detector is employed to locate the feature points in the frames. Secondly, the matched key-points are used for estimation of affine transform to reduce false matches. Finally, motion estimation is performed based on the affine model and the compensation for vibration is conducted based on spline smoothing. It was reported that this algorithm can process up to 30fps on a Workstation with an Intel Xeon 2.26 GHz processor and 6 GB RAM for images with resolution of 320x240 pixels.

A very recent work by Dong et al (2014) is able to process a 640x480 pixels video at a speed of 40fps on a notebook computer with a 2.5GHz Intel Duo Core CPU. While it is able to achieve a high processing speed, the improvement comes at the price of unreliable motion estimation, which leads to failure in stabilizing the video ultimately. Moreover, the algorithm proposed does not furnish evidence for real-time processing on real-time on-board systems.

It is noted that these previous methods work offline or in a post-process way. For some UAV vision tasks, frames are required to be stabilized immediately and be presented to other tasks, such as object tracking and detection.

The proposed method in this paper is closely related to the method proposed by Ho (Ho (2014)). The algorithm employs finding corners in frames, followed by estimating a 2D motion model between consecutive frames. Then the parameters in motion model are treated as trajectory to be smoothed using an averaging window. This smoothing of the motion model parameters is different from the motion trajectory smoothing employed by previous works (Shen et al (2009); Vazquez and Chang (2009); Wang et al (2011); Dong et al (2014); Lucas et al (1981)). Our framework is similar to Nghia Ho's but has a significant difference in terms of implementation and performance. Most notably, a novel hybrid mechanism for motion estimation and an optical flow-based corner tracker has been proposed to overcome the challenges encountered by previous algorithms. In addition, the proposed stabilization algorithm performs in real-time, whereas Nghia Ho's algorithm processes the whole video sequence before achieving stabilization.

The rest of this paper is structured as follows. In section 2, we introduce the proposed framework for video stabilization. Then in section 3, we discuss how the processing speed of motion estimation is expedited by using an optical flow-based corner tracker and a reduced region of interest. Then in section 4, we explain the improvisation of the motion estimation given by optical flow-based tracker with hybrid motion estimation mechanism. Following which, we verify the speed of the algorithm on publicly available UAV videos and achieve a speed faster than the current state-of-the-art algorithm. In section 5, we discuss the implementation of the real-time component of the algorithm using multi-thread processing. Finally, we conclude the paper in section 6.

---

## 2 Proposed Framework

Spurred by the challenge to meet the real-time requirements, we propose a novel stabilization framework as shown in Figure 2. The idea is to build an appropriate model for the global motion trajectory of the UAV camera, estimate the motion parameters between two successive frames using very efficient key point detector and compensate the motion via an effective optical flow based tracking of the key points.

Firstly, the key novelty of the proposed algorithm lies in the high speed estimation of the motion parameter trajectory between two consecutive frames, which other methods are less efficient in accomplishing. This is because the motion estimation phase chooses an appropriate lower order homography, either rigid or similarity transformation between frames. This is so that the algorithm is able to frequently use lower order of homography to estimate the motion parameters unless there is a need to use higher order. As a result, achieving a high processing speed. Secondly, our motion compensation phase smooths the parameters via an averaging window. Lastly, we warp the frame according to the smoothed parameters. The program is able to process up to 100 frames per second for all 3 stages of the algorithm. This reassures us that the algorithm is suited for real-time video stabilization.

After achieving such a high speed in stabilizing the frames, the real-time version of the algorithm can be implemented. The three stages of video stabilization were separated into three different threads. Threads are basically independent processes that can run concurrently and share the same resources inside a program. Our algorithm is implemented such that motion estimation thread, motion compensation thread and image composition thread are synchronized in such a way that they are able to stabilize a UAV video in real time. In other words, there is no need to wait for the processing for the whole video to be done.

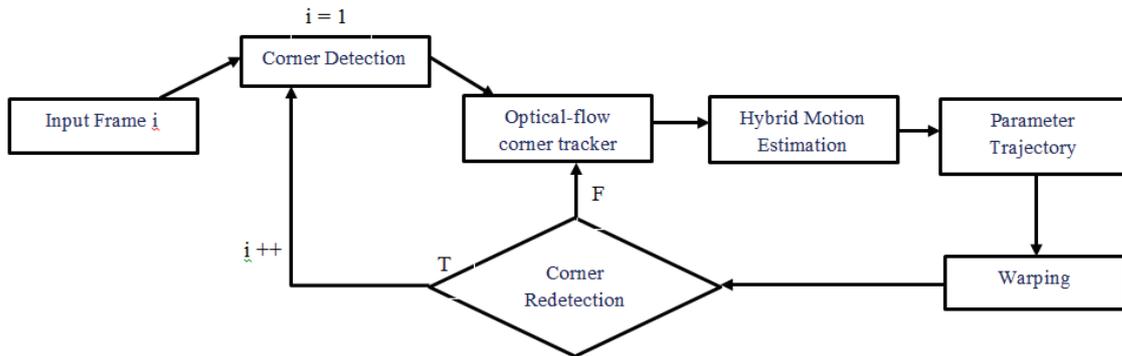


Fig. 2: Overview of the proposed video stabilization framework

### 3 Optical-Flow-Based Tracking

Motion estimation is the very first step in video stabilization and it is also the most time consuming step. Consequently, our primary goal is to cut down the time required to calculate the motion parameters trajectory between consecutive frames in the video. The motion trajectory is estimated from how the feature points (Harris and Stephens (1988)) move between consecutive frames. Therefore, to initialize motion estimation, we first efficiently detect feature points.

#### 3.1 Feature Point Detection

Feature points are locations in the image with large variations in intensity in all directions. They are very important in motion estimation step as they determine the quality of the motion estimation. One early attempt to find these corners was done by Chris Harris and Mike Stephens in their paper (Harris and Stephens (1988)), which now is called Harris Corner Detector. It basically finds the difference in intensity for a displacement of  $(u, v)$  in all directions. This can be easily expressed as below.

$$E(u, v) = \sum_{x,y} w(x, y) [I(x + u, y + v) - I(x, y)]^2 \quad (1)$$

where  $w(x, y)$  is either a rectangular window or Gaussian window function which gives weight to the surrounding pixels. In 1994, Shi and Tomasi (1994) made a small modification to it in their paper which shows better results compared to Harris Corner Detector. The scoring function in Shi-Tomasi Corner Detector is given by:

$$R = \min(\lambda_1, \lambda_2) \quad (2)$$

- When  $|R|$  is small, which happens when  $\lambda_1$  and  $\lambda_2$  are small, the region is flat.

- When  $R$ , which happens when  $\lambda_1 \gg \lambda_2$  vice versa, the region is edge.
- When  $R$  is large, which happens when  $\lambda_1$  and  $\lambda_2$  are large and  $\lambda_1 \approx \lambda_2$ , the region is a corner.

Since the proposed idea is to ensure that the algorithm is able to run in real-time, a simple experiment of mosaic was conducted. Its purpose is to find out if the number of key points detected has a huge effect on the estimated transformation parameters and also to find out how the time is affected when the number of key points returned is scaled up.

In this experiment, two images of the same scene will be stitched together with a reference image as shown in Figure 3.

Firstly, key points will be detected in the left image and the reference image. Next, key point matching will be performed between the left image and the reference image. Finally, a homography will be generated to transform image 1 to the perspective of the reference image. The cycle repeats itself for the right image.

If observed carefully, there is not a substantial qualitative difference between a 200 Key Points Mosaic image and a 1000 Key Points Mosaic image as illustrated in Figure 4. The quality of both mosaics are similar. This demonstrates that large number of Harris corners is unnecessary. Therefore, keeping in mind that our proposed idea is to run the algorithm in real-time, the number of key points needed for a good estimate of the motion model is set to be less than 200. From experiments, we set a value of 50 which gave us visually pleasing results without compromising on stabilization quality.

Next, the area for detection of points is implemented such that the key point detection will be from an area smaller than that of the frame (see Figure 5). The region of interest will be set such that the length and width of the region are smaller than the size of the frame. The reason why this region is implemented is because the pixels close to the edges of the frame have a high probability of not appearing in the next consecutive frame, especially when the UAV is in motion. In this way, it can prevent

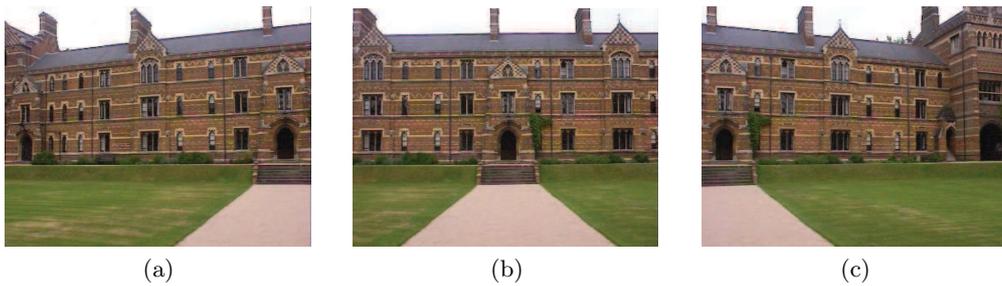


Fig. 3: Mosaic Images



Fig. 4: Mosaic Results



Fig. 5: Reduced Region of Interest

the key points at the edges from being detected, saving computational time needed to match those key points which have a high probability of vanishing in the next frame. The next step is to match the key points.

### 3.2 Feature Point Matching

To speed up the matching, a common way is to use Approximate Nearest Neighbors (Andoni and Indyk (2006)). This method is still slow for our real-time application. Furthermore, these techniques are only useful when the objective is just to characterize the neighbourhood rather than the exact locations. In image stabilization, locations are important so that the homography generated is accurate. Therefore, we propose an optical flow based matching of corners detected in one frame to the next frame.

Optical flow (Fleet and Weiss (2006)) is the pattern of apparent motion of image objects between two con-

secutive frames caused by the movement of an object or camera. It is a 2D vector field where each vector is a displacement vector showing the motion of points from first frame to second. Optical flow assumes brightness constancy, as shown below.

$$I(x, y, t) = I(x + dx, y + dy, t + dt) \quad (3)$$

Then taking Taylor series approximation of right-hand side, and removing common terms and dividing by  $dt$  to get the following equation:

$$f_x u + f_y v + f_t = 0 \quad (4)$$

where:

$$f_x = \frac{\partial f}{\partial x}; f_y = \frac{\partial f}{\partial y} \quad (5)$$

$$u = \frac{dx}{dt}; v = \frac{dy}{dt} \quad (6)$$

The above equations describe the optical flow in terms of the spatial image gradient. In it, we can find  $f_x$  and  $f_y$ , they are image gradients. Similarly,  $f_t$  is the gradient along time, but  $(u, v)$  is unknown. We cannot solve this one equation with two unknown variables. Various methods have been suggested to resolve this problem and we choose the gold standard, Lucas-Kanade algorithm (Lucas et al (1981)).

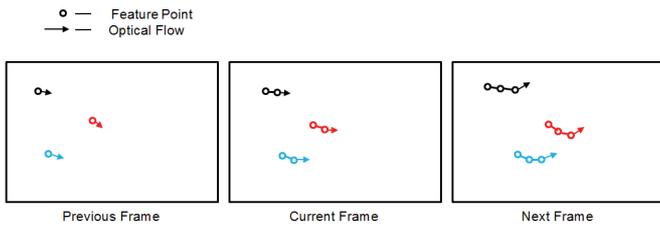


Fig. 6: Optical-Flow to estimate motion

By estimating the movement of the corner points, we do not have to detect corners in the next frame. And since we have estimated where the corners are in the next frame, matching is already completed for us. Optical flow matching is repeated for five frames with corner detection only for the first frame, and then the corners are re-detected in the sixth frame. This is done to not violate the two main assumptions of optical flow, which are:

1. The pixel intensities of an object do not change between consecutive frames.
2. Neighbouring pixels have similar motion.

Using the Lucas-Kanade method, we obtain the optical flow vectors of the corner points in the first frame and subsequently for the next four frames. But again, there are some problems. Until now, we were dealing with small motions and the assumptions fail when there is large motion. So we use pyramids (Adelson et al (1984); Bouguet (2001)) to remove small motions and large motions becomes small motions as we go higher in the pyramid. Now applying Lucas-Kanade, we get the optical flow along with the scale.

We detect the Shi-Tomasi corner points in the first frame. Then, we track those points using Lucas-Kanade optical-flow in iteration. In subsequent frames, the tracker will provide us with the new locations of the corner points. We pass these next points as the previous points in the next step.

The motion vectors are the coordinates of the key points that existed in the previous frame and have moved to a new location in the current frame. After retrieving the new key point locations, we know that there exists a set of corners that are both in the previous frame and the current frame. This permits us to generate an inter-frame transformation matrix that is used for motion modeling and image composition.

## 4 Homography Estimation

After acquiring the corners that exist in both current and previous frames, we can now estimate the motion model between these frames. As mentioned above, we know the pixels in the previous frame have moved in the current frame. Using this information, the motion estimator takes in the coordinates of the feature points

in both frames to generate a homography. This function allows the homography returned to be of similarity (also known as partial affine) or rigid (also known as Euclidean) transformation and it is implemented such that the homography returned depends on the hybrid motion estimation used. The deciding factor of the mechanism will be discussed in the next following subsection. If the homography returned is a rigid transformation, the transform parameters will store only the  $t_x$ ,  $t_y$ , angle of the transformation and a scale factor of 1. If the homography matrix returned is of similarity transformation, the transform parameters will store  $t_x$ ,  $t_y$ , angle and an arbitrary scale factor of the transformation. After we extract these transform values, we treat them as motion parameters.

### 4.1 Hybrid mechanism for motion estimation

In previous works, motion estimation is carried out either using particle filters (Yang et al (2006)), or variational methods (Pilu (2004)), and mostly using SIFT feature based matching techniques (Battiatto et al (2007)), all of which are not suitable for the real-time implementation considered in this work. An extensive discussion of this point can be found in (Dong et al (2014)). In fact, Dong, pointed out that the requirement of the long-range feature tracking makes many existing methods incompetent for challenging cases, since long feature trajectories are difficult to obtain in sequences with rapid scene changes, texture less objects, severe occlusions, or excessive motion blur. All these scenarios are common for UAV videos and they are tackled in this work. Citing the above reasons, Dong made use of a fast KLT tracker and achieved an average speed of 50 fps on publicly available UAV videos.

In general, the motion estimation finds an affine transform  $[A-t]$  (a  $2 \times 3$  floating-point matrix) that approximates best the affine transformation between two sets of points. In case of feature point sets, the problem is formulated as follows: you need to find a  $2 \times 2$  matrix  $A$  and  $2 \times 1$  vector  $t$  for the source points (src) and the destination points (dst).

$$[A^* | t^*] = \arg \min \sum_i ||dst[i] - Asrc[i]^T - t||^2 \quad (7)$$

Solving for  $[A-t]$  requires a minimum of 3 pairing points that are not degenerate. This is straight forward to do. Lets denote the src point to be  $X = [x \ y \ 1]$  and the dst to be  $Y = [x' \ y' \ 1]$ , giving:

$$TX = Y$$

$$\begin{bmatrix} a_{11} & a_{12} & t_x \\ a_{21} & a_{22} & t_y \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} \quad (8)$$

We can re-write this as a typical  $Az = b$  matrix and solve for  $z$ . We'll also need to introduce 2 extra pair of points to be able to solve for the motion parameters, as follows.

The affine transform mentioned earlier has a degree of freedom of six. There exist two lower degree of freedom transformations, namely similarity and rigid transformation. Those two transformations have four and three degrees of freedom respectively. The four degrees of freedom include rotation angle, scaling, translation in  $x$  and translation in  $y$ . Rigid transformation assumes no scaling.

Figure 7 illustrates the idea of the hybrid mechanism.  $E_r$  and  $E_s$  represent the total root-mean-square distance, ( $e_r, e_s$ ), between the  $N$  corners of the previous frame and current frame, after applying the respective transformations. These two values will be computed and compared. If the root-mean-square distance for rigid transformation is lower than similarity transformation, then the mechanism will select rigid transformation computation for a sequence of frames (i.e. dwell time). To achieve a balance between speed and stability, a dwell time of 20 frames is used to switch between rigid (3 parameters) and similarity transformation (4 parameters).

For every 20 frames, when this mechanism is used, key points that are detected will be used together with the two homography matrices, namely rigid transformation matrix and similarity transformation matrix. At the initial frame both rigid and affine matrices are estimated from the set of key points that are detected and matched from optical flow. Then, we estimate two sets of new motion vectors using the two transformation matrices. The next step in the algorithm will then calculate the average distance between the new points and the key point locations for both the transformations. If the average distance between the points generated for rigid transformation is smaller than that of the average distance between the points generated for similarity transformation, then the mechanism will remain to be rigid transformation, or switch to similarity transformation. Therefore, the algorithm will return the transformation that gives lesser change in distance between the points.

The hybrid mechanism is switching between rigid and partial affine transformation. This means that we are stabilizing shaky motions using only up to 4 degrees of freedom, namely, translation in  $X$  direction, translation in  $Y$  direction, the scale factor and the rotation factor.

Based on the video stabilization algorithm described so far, we have accomplished the implementation of the homography hybrid mechanism. The hybrid mechanism has achieved two main objectives:

1. Reduce movement of the features, thereby reducing jitters between two successive frames.
2. Decrease computational time by choosing the transformation with less parameters whenever possible.

Figure 8 gives an overview of motion estimation work flow in the proposed algorithm. It begins by detecting

corners in frames, then utilizes optical-flow to detect and track corners in subsequent frames. On a rare basis, if the number of corners detected is too less, then stabilization is skipped for that frame. The optical flow tracking is followed by weeding step which checks for flow consistency in the backward direction, thereby eliminating bad matches. For every twenty frames, hybrid mechanism is called upon to determine the best transformation to be estimated. Then the parameter trajectories are extracted from the estimation and stored. The process is repeated for the number of frames in the video.

---

## 5 Motion Compensation and Image Composition

### 5.1 Averaging Window Smoothing

In the motion compensation stage, we employ a basic method of accumulating trajectory parameters and then smoothing them using an averaging window. The trajectory of the parameters is accumulated within a window of frames, which is double the size of smoothing radius, and then the result of the addition of each transformation parameters is then averaged out according to the window size. The result of the averaging of the parameters is then the smoothed trajectory. Next, iterate through each of the original trajectory and adding the difference between the smoothed trajectory and the accumulated global trajectory.

Figure (a) of 9 shows the values of  $t_x$ . It can be viewed in Figure (b) of 9 that the values of  $t_x$  fluctuate a lot, which correlates with our observation of the video being shaky. The next step is to smooth the  $x$  values separately after accumulating the frame-to-frame transformation using a sliding average window. It is worth mentioning that the parameter trajectory is a rather abstract quantity that doesn't necessarily have a direct relationship to the motion induced by the camera. For a simple panning scene with static objects, it has a direct relationship with the absolute position of the image. The important insight is that the trajectory can be smoothed, even if it does not have any physical interpretation. The result of the smoothing process is shown in Figure 9.

In the final step of the video stabilization algorithm, a basic method is used to couple warping with cropping. First, warp the frame according to the new set of homography derived from the previous step. Next, the cropping and resizing will reduce unfilled area cause by the smooth motion parameters.

---

## 6 Real-Time Implementation

In the previous sections, it is discussed that the algorithm works by these 3 steps, motion estimation, motion compensation and image composition. These 3 steps are usually executed one after another. This means that the

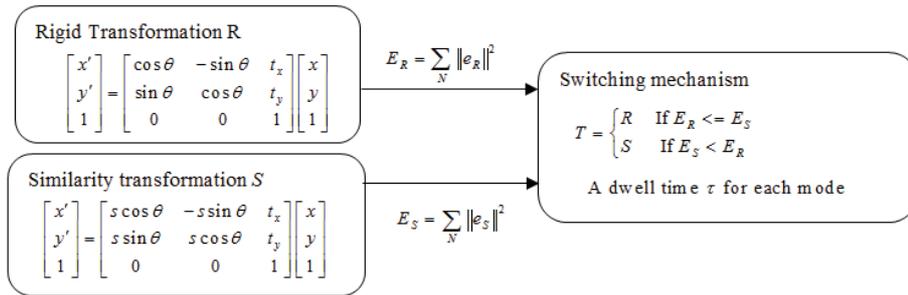


Fig. 7: Hybrid Mechanism for Motion Estimation

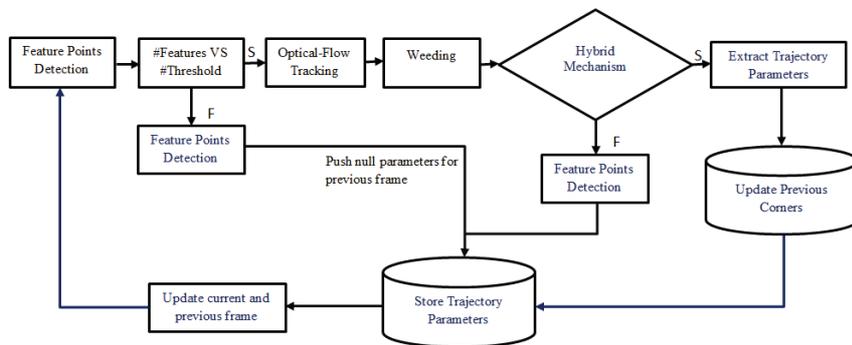


Fig. 8: Motion Estimation

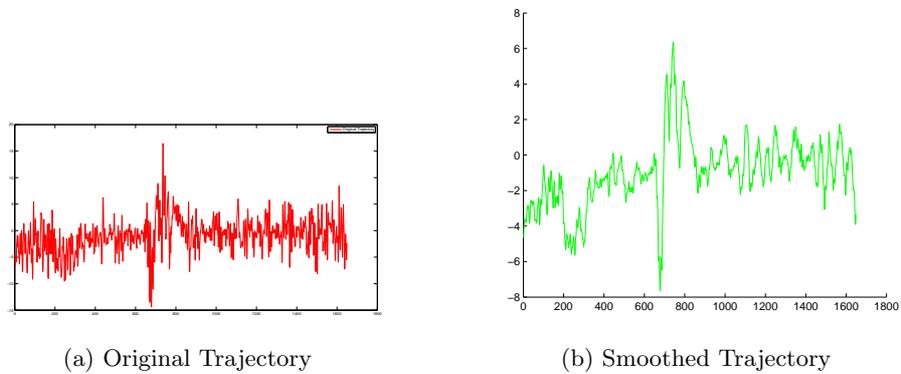


Fig. 9: Compensation Results

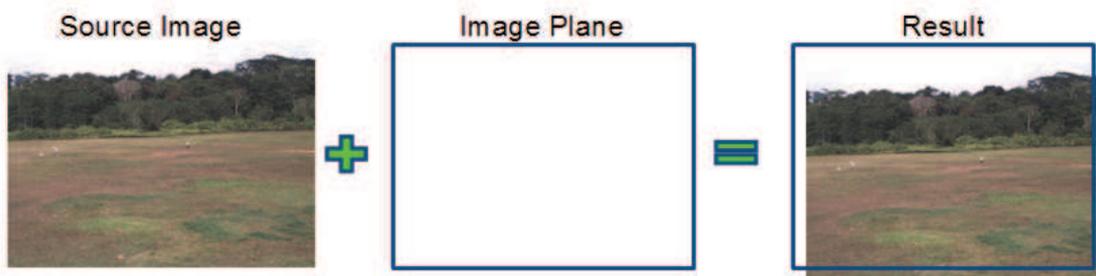


Fig. 10: Image Warping into Image Plane

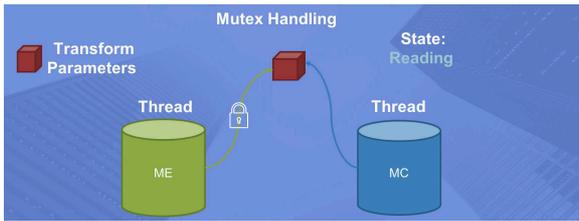


Fig. 12: Multi-Thread Mutex Control

algorithm has to have the entire video as the input, then it will process all the frames within the video via those steps mentioned and then display the output, which is not real-time at all.

This brings us to the second phase of this development which is to make the algorithm real-time as shown in Figure 11.

### 6.1 Multi-Threaded Approach

The program consists of three threads as shown in Figure 12, which are capable of running independently of each other. Parallel computing is used to separate the video stabilization algorithm into three parts: thread 1 (motion estimation), thread 2 (motion compensation) and thread 3 (warping). Using the first thread, the motion estimation is carried out for the first twenty incoming frames (smoothing radius for motion compensation is set to ten) while the other threads wait for this process to complete. This delay is because the motion compensation and warping require smoothed corner trajectories that can be generated only after the first twenty frames are processed. From the 21<sup>st</sup> frame, all the three threads operate simultaneously. In other words, the motion estimation between the 20<sup>th</sup> and 21<sup>st</sup> frame is obtained using thread 1 and the other threads compute the smoothed trajectory using the information generated from the second frame to the 21st frame. The whole process continues indefinitely until the video stream ends. Since the threads themselves process the frames much faster than the frame rate, the timings of the threads eventually catch up with the timings of the frame availability.

### 6.2 Motion Estimation Thread

Once the ME thread is able to open a video file, it will keep estimating the homography matrix between each frame until it reaches the end of the video file. The motion compensation thread will be notified to run as soon as ME thread finished estimating for twenty frames.

### 6.3 Motion Compensation Thread

The motion compensation thread averages the transform parameters within the storage for the transform param-

eters. From now on, whenever a new frame is estimated by ME thread, the motion compensation thread will remove the oldest transform parameters, insert the latest parameters and average. It is to be reminded that the averaging at motion compensation stage has not changed from the previous implementation. The stage no longer has the whole videos motion trajectory parameters to work with anymore. Instead, it collects and manages the data storage of the transform parameters given by motion compensation thread.

### 6.4 Image Composition Thread

Once motion compensation produces the smoothed parameters of the latest frame, it will insert the data into the storage so that the image composition thread is able to retrieve the data and frame to display them accordingly.

## 7 Results

### 7.1 Results of Stage I

We tested the enhanced corner tracking and hybrid motion estimation system on the video obtained from in-house UAV videos and on standard benchmark videos (Dong et al (2014)). Figure 13 shows a screenshot of the comparison between the unstable and the stabilized video on the Forest video from the public database.

A web demo for the Forest video is shared below:

<http://tinyurl.com/jv6fvvu>

Compare the above result to the result obtained using rigid and similarity transformation alone:

<http://tinyurl.com/gnploac> (rigid)

<http://tinyurl.com/glxt65m> (similarity or partial affine)

It is evident that the hybrid motion estimation is much better in comparison to the simple rigid transformation or similarity transformation. This can be clearly witnessed after 20 seconds into the video when there are sudden camera movements.

More tests were conducted on UAV videos available as public video stabilization datasets. These tests were conducted on a desktop with a CPU clock speed of 1.7GHz. The videos used are aerial view videos taken from a UAV to simulate environment when a real UAV is used for surveillance. This algorithm is able to estimate the motion matrix of adjacent frames at a maximum speed of 100 frames per second. The same algorithm is used to estimate the motion matrix for a frame size of 640x480 at a speed of 30 frames per second. This kind of result reassures us that this algorithm is very much suited for real-time processing of motion. Below are the results of how many times rigid and affine transformations have been utilized in different videos.

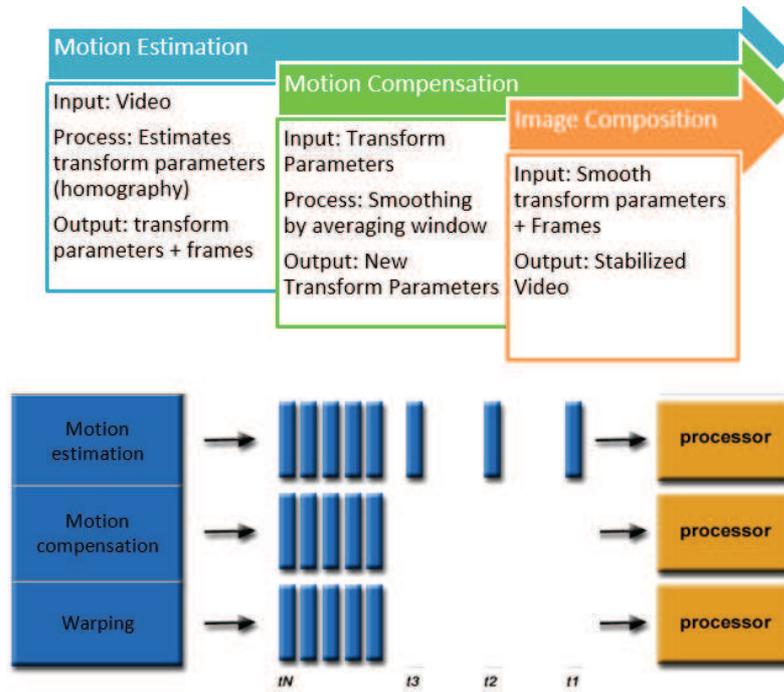


Fig. 11: Multi-Thread Concurrent Work Flow



Fig. 13: Before and After Stabilization

Video	Number of Rigid Transformations (per 20 frames)	Number of Similarity Transformation (per 20 frames)
Forest	41	41
Ground	16	32
Bird	21	9
Girl	12	10

Table 1: Number of times different motion models are used

It is clear that the hybrid motion estimation uses both rigid and similarity transformations in a manner dependent on the nature of the motion. For motion behaviour that is not easy to be captured by the simple rigid mechanism, the similarity transformation offers a better model for motion behaviour. The web-demos of the other videos in the public database can be found below.

1. Ground: <http://tinyurl.com/zfb4lo6>

2. Bird: <http://tinyurl.com/jsn7g24>

3. Girl: <http://tinyurl.com/hodl9ey>

4. Hippo: <http://tinyurl.com/zrq8jna>

Using video frames obtained from in-house UAV videos, we tested the hybrid motion estimation based stabilization algorithm. The captured video frames were obtained by a UAV that shows the motion of another UAV. The video frames suffer from both object motion and scene

motion (due to the on-board cam motion). Below is a demo of our system.

<http://tinyurl.com/z8muwab>

From the above results, we can see that after the first few seconds, our method can better stabilize the unstable frames, especially when seen from the point of view of the UAV in the picture.

Furthermore, we observe that although the scene changes abruptly, our system can handle them without crashing. This is a common scenario when mechanical gimbals do not provide much stability to the captured UAV videos or when there are frame rate issues or when there is a sudden mid-flight recourse taken by the UAV. All these scenarios were handled efficiently by our system, as demonstrated in the comparison video above.

## 7.2 Motion Estimation Timings

In order to ensure that our algorithm has the potential to execute in real-time, we time the duration needed to finish estimating the motion model for the publicly available videos. Most of the algorithms reviewed in the literature are either not applicable for real-time processing or needs more processing time than the proposed algorithm in this paper. Table 2 compares the processing speed of the proposed algorithm and the latest instantaneous video stabilization method in the literature. For a fair comparison, we use a notebook computer with a dual-core 1.70 GHz processor, as done in Dong et al (2014), using a notebook computer with a 2.5 GHz Intel dual-core CPU. Motion estimation step is the most expensive step in video stabilization. We are able to utilize optical-flow-based tracking to reduce the computational time.

In Table 2, we recorded the timing for the motion estimation stage to finish computing and smoothing the global motion trajectory of the video. From the table, we can infer that the lower the resolution of the video the faster it is to calculate the motion model in between the frames. Furthermore, we speed up the algorithm by using tracking which effectively remove the need to do iterative detection of corners for each frame. This is how we achieve such an impressive timing for motion estimation.

## 7.3 Results of Real-Time Implementation

Multi-threaded implementation is possible due to the high speed of the proposed stabilization framework compared to existing works. Experiments were conducted to estimate the speed at each stage. The result was that motion compensation and image composition step combined takes much shorter time than motion estimation (ME). Hence it is concluded that ME is the bottleneck of the operation. Since through experimental results, ME can reach up to 81 frames per second on a 640x480p video,

this implementation can allow the program to process frames in real-time.

However, it has to be kept in mind that the threads execute under the same program explained in section 6. This means that they are able to share resources or data with one another. This feature is perfect for the proposed algorithm because motion compensation needs to make use of the output (transform parameters) from the ME process and image composition needs to make use of the output from motion compensation as seen from Figure 11.

This feature can serve as a convenient way to read and write to the same resource in the memory space or it can be dangerous as data corruption can happen if the data accessibility is not handled carefully. In this program, mutex is used to facilitate resource accessibility. A demo of the in-house UAV can be viewed using the link below.

<http://tinyurl.com/hel53sb>

Furthermore, with the usage of multi-threading processing, the time involved to compute and smooth the parameter trajectory is reduced significantly, as shown in Table 3. This is because once the parameters are calculated, the motion compensation thread is able to smooth the parameter values immediately without having to wait for the computation of the parameter trajectory of the entire video to be extracted.

This section contains the final video results obtained from our algorithm which are available online.

1. Ground: <http://tinyurl.com/zfb4lo6>
2. Bird: <http://tinyurl.com/jsn7g24>
3. Girl: <http://tinyurl.com/hodl9ey>
4. Hippo: <http://tinyurl.com/zrq8jna>

---

## 8 Conclusion

This paper describes the algorithms and methods used to tackle the problem of real-time video stabilization for unmanned aerial vehicles (UAVs) videos. Due to the size and structure limitations, UAVs are highly susceptible to atmospheric turbulence, which induces jitters and makes the video unstable. It is very difficult to detect and track targets of interest in such unstable videos. Therefore, it is necessary to design a real-time video stabilization algorithm for UAVs as a pre-processing module for higher-level vision tasks, such as object detection and tracking. To achieve the above goal, two main components were designed as part of the proposed video stabilization algorithm: (1) By designing an appropriate motion model for the global motion of the UAVs, the proposed stabilization mechanism avoids the necessity of estimating the most general motion model, projective transformation, and considers simpler motion models like the similarity and rigid transformation; (2) In order to achieve high processing speeds, an optical flow based motion estimation is proposed to replace the conventional tracking and

Video Name	Resolution	Frames	Dong et al (2014)(fps)	Offline Method(fps)
Forest	320x240	1650	77.2	107.8
Ground	640x480	966	39.9	34.9
Bird	640x360	601	47.4	44.2
Girl	640x360	446	42.7	45.5

Table 2: Time taken to complete pre-process entire video

Video Name	Resolution	Frames	Offline Method(fps)	Real-Time Version
Forest	320x240	1650	107.8	126.32
Ground	640x480	966	34.9	58.93
Bird	640x360	601	44.2	72.65
Girl	640x360	446	45.5	65.36

Table 3: Comparison between offline and real-time

matching algorithms used by state-of-the-art video stabilization methods. These novel ideas resulted in a real-time stabilization algorithm. A demo of the stabilized videos can be accessed using the web links provided in this paper.

Our method will not yield excellent results for videos that has large and gradual scene changes. This will cause the frames to be warped in such a way to resist the changes resulting in large unfilled areas in the frames. Future work will be carried out to resolve this issue.

## References

- Adelson EH, Anderson CH, Bergen JR, Burt PJ, Ogden JM (1984) Pyramid methods in image processing. *RCA engineer* 29(6):33–41
- Andoni A, Indyk P (2006) Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions. In: 2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06), IEEE, pp 459–468
- Battiatto S, Gallo G, Puglisi G, Scellato S (2007) Sift features tracking for video stabilization. In: *Image Analysis and Processing, 2007. ICIAP 2007. 14th International Conference on*, IEEE, pp 825–830
- Bouguet JY (2001) Pyramidal implementation of the affine lucas kanade feature tracker description of the algorithm. *Intel Corporation* 5(1-10):4
- Dong J, Xia Y, Yu Q, Su A, Hou W (2014) Instantaneous video stabilization for unmanned aerial vehicles. *Journal of Electronic Imaging* 23(1):013,002–013,002
- Fleet D, Weiss Y (2006) Optical flow estimation. In: *Handbook of mathematical models in computer vision*, Springer, pp 237–257
- Harris C, Stephens M (1988) A combined corner and edge detector. In: *Alvey vision conference*, Citeseer, vol 15, p 50
- Ho N (2014) Simple video stabilization using opencv. URL <http://nghiaho.com/?p=2093>
- Lucas BD, Kanade T, et al (1981) An iterative image registration technique with an application to stereo vision. In: *IJCAI*, vol 81, pp 674–679
- Pilu M (2004) Video stabilization as a variational problem and numerical solution with the viterbi method. In: *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*
- Shen H, Pan Q, Cheng Y, Yu Y (2009) Fast video stabilization algorithm for uav. In: *Intelligent Computing and Intelligent Systems, 2009. ICIS 2009. IEEE International Conference on*, IEEE, vol 4, pp 542–546
- Shi J, Tomasi C (1994) Good features to track. In: *Computer Vision and Pattern Recognition, 1994. Proceedings CVPR'94., 1994 IEEE Computer Society Conference on*, IEEE, pp 593–600
- Vazquez M, Chang C (2009) Real-time video smoothing for small rc helicopters. In: *Systems, Man and Cybernetics, 2009. SMC 2009. IEEE International Conference on*, IEEE, pp 4019–4024
- Wang Y, Hou Z, Leman K, Chang R (2011) Real-time video stabilization for unmanned aerial vehicles. In: *MVA*, pp 336–339
- Yang J, Schonfeld D, Chen C, Mohamed M (2006) Online video stabilization based on particle filters. In: *2006 International Conference on Image Processing*, IEEE, pp 1545–1548