

C-14: Assured Timestamps for Drone Videos

ABSTRACT

Inexpensive and highly capable unmanned aerial vehicles (aka drones) have enabled people to contribute high-quality, aerial videos at a global scale. However, for videos to be trusted we must verify that they are 'real' videos rather than synthetic or modified videos. But establishing that videos are real is not enough. A key challenge exists for accepting videos from untrusted sources: establishing *when* a particular video was taken. Once a video has been received or posted publicly, it is evident that the video must have been created before that time, but there are no current methods for establishing how old it is.

We propose C-14,¹ a system for instructing a drone to take videos in a particular pattern that helps establish the earliest time the video could have been created. Using recent developments in optical flow and camera pose estimation, we show that we can establish such a time using both unconstrained flight patterns with many waypoints and drone motions, as well as constrained patterns that employ steady, easy to verify motions. Through extensive experiments we show how to verify a 59-second unconstrained video with eight motions in 91 seconds of computation with a false positive rate of one in 10^{13} and no false negatives. We also verify a 190-second constrained video with 4 motions in 158 seconds of computation with a false positive rate of one in one hundred thousand and no false negatives.

1. INTRODUCTION

The continuous proliferation of drone-mounted, high resolution video cameras is ushering in an era of *global scale video sensing*. For instance, drones have enabled citizens to provide video coverage of land areas that were previously inaccessible. In areas of Latin America and Asia, citizen-powered drones are already proving crucial in providing imagery and video to monitor large areas of land vulnerable to unauthorized development such as deforestation [43]. Similarly, freelance journalists are producing invaluable evidence from war zones and other difficult to cover areas from the air [1].

However, a key challenge is ensuring the trustworthiness

of videos. We divide this problem into three parts: (i) location establishment: the video was taken at a particular place, (ii) integrity assurance: the video was unaltered after recording (e.g. spliced, re-encoded, etc.), and (iii) timestamping: a video was taken at a particular time. A number of research efforts address the first [23,37,47] and second problems [25,27,35,45], but we are unaware of techniques that address the third: *when* a video was taken. For example, while nearby infrastructure can attest to a location [37], the drone may provide a very old video instead of a fresh one.

We propose a system, C-14, that assures a video was not created before some time t_b or after some time t_e . Demonstrating t_e is straightforward using any trusted storage mechanism, such as a public blockchain ledger. The video creator can simply post a hash of the video to a blockchain at t_e , thus proving that the video was created at some time previous. Providing proof that the video was taken after t_b is much more difficult. If one posts a video's hash at t_b , the video creator could have created the video a long time ago and simply held it until t_b .

To establish that a video was taken after a time t_b , C-14 introduces a pattern of movements in the video that could only be known by the creator after time t_b . This pattern of movements can be provided by a trusted party only after t_b or could be derived from values publicly known after t_b . For instance the hash from a blockchain could only have been known at the time the block was produced or after. From that value we derive a motion pattern and prescribe that the drone video matches that motion pattern. This demonstrates that the video could only have been produced after t_b .

The pattern is a series of translations (move up/down, left/right, forward/backwards), or rotations (yaw, roll, tilt), and combinations of the two. We demonstrate the system on two types of patterns: highly variable drone motions over wide areas, and highly focused and controlled motions over a small area. The advantage of highly variable motions is that verification can be incorporated into typical drone flights, however controlled motions are more amenable to sampling small parts of the video, which speeds verification.

¹The name C-14 comes from the radioactive isotope used for carbon dating organic matter.

C-14 incorporates a number of techniques to speed verification: compression, skipping frames, spatial sampling, and temporal sampling. A full analysis of the video to verify the motion runs in 2000x real-time (a 2 minute video takes 3 days of GPU time to verify), but through compression and sampling we can reduce the amount of time needed. C-14 can verify a 59-second unconstrained video with 8 motions in 91 seconds of computation with a false positive rate of 1 in 10^{13} and a 190-second constrained video with 4 motions in 158 seconds of computation with a false positive rate of 1 in one hundred thousand and no false negatives.

2. ASSUMPTIONS AND THREAT MODEL

We have made certain assumptions in building C-14:

- **No Trusted Hardware:** C-14 does not rely on any trusted hardware on the drone, such as a TPM. Such hardware is not yet readily available on drones and it would need to encompass the clock, flight controller, and video sensor to be effective at establishing a timestamp. Instead we have targeted off-the-shelf devices.
- **Flight Pattern Unknown Before t_b :** The flight pattern must have sufficient entropy as to make it unguessable. The pattern can be derived from an online public source, such as a blockchain, or from a trusted third party that reveals a random seed or full flight pattern at a certain time.
- **Limited Motion:** The system has been built and tuned around largely static scenes, such as buildings, forests, etc. Drone flights in the USA cannot be done over live subjects, which makes quantifying this limitation difficult. However, techniques to remove independent object motion from scenes [12,13,20,41] could be used to increase the robustness of C-14.
- **Featured Scenes:** We assume that the subject of the video is not mono-tone or mono-textured, such as a field of snow, a body of water, or a featureless desert. Such scenes pose difficulties for the optical flow algorithms we use in C-14.
- **Speed-Agnostic:** We do not make any assumptions about the speed of the drone. Verification of the video is done on a frame basis, not based on real time.

We assume a particular threat model to the system, and classify many attacks as solvable through other mechanisms. Most importantly C-14 can only assure that the video was created after t_b . It cannot, on its own, assure the integrity of the video, or its location or subject matter. Intuitively, altering or forging videos is more difficult than forging static images, and a great deal of research has shown the creating fraudulent images that stand up to scrutiny is incredibly difficult. We detail a few of these attacks, but the defenses are outside of the scope of this paper.

- **Video Splicing:** An attacker could take a large number of videos from a location, each of which performs one of the movements required from the authenticated video. Once the sequence of movements is known, one can splice together those movements into a video and present it as authentic. However, such splices can be detected in videos through a number of techniques [9–11,21,22].
- **3D Rendering/Fake Video:** If an attacker can generate a video that is a full reconstruction of a scene, they can arbitrarily create any pattern (rotations and translations). Detecting videos created from whole-cloth is outside of the scope of our work, but the image forensics community has worked on detecting such reconstructions based on a number of techniques, for example camera noise [17], the smoothness of images [30] and machine learning classification [31]. A similar technique would encompass such videos created by so-called “bullet-time” or 3D reconstructions from large numbers of photographs.

Most importantly it is best to think of C-14 in the same light as a CAPTCHA—we can provide some assurance and raise the bar for an attacker, but any such system will be a continuous game of cat and mouse if the stakes are sufficiently high. C-14 is one piece of a system for assuring properties of drone videos.

3. DRONE BACKGROUND

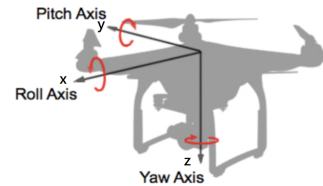


Figure 1: This figure demonstrates the axes of movement found in a copter drone. Image from DJI documentation [2]

Here we provide background on how drones fly and collect videos. We focus on ‘copter’ drones which typically have four rotors that allow the drone to translate in three dimensions, as well as rotate around three axes, as labeled in Figure 1. The drone is equipped with a front-facing video camera, mounted on a three-axis gimbal. The gimbal can be manually pointed in any direction, but typically it is set to *yaw-follow* meaning that the gimbal tries to maintain a constant pitch angle and zero rotation with respect to the world horizon. The yaw of the camera follows the drone’s heading, though it does so with some elasticity to prevent sudden motions in the video. Consider what happens when the drone

moves to the right (a positive y -axis translation). The aircraft adjusts the speed of the propellers to roll slightly right (a positive roll), which makes the aircraft move to the right. The gimbal counteracts this motion and the resulting video has no roll.

For the remainder of the paper we use a frame of reference we refer to as the *ideal drone* frame of reference. This frame of reference is the drone without roll and pitch induced by aircraft motion. In other words the drone would always appear to have no roll and pitch in this frame of reference, but it does yaw.

4. MOTION PROGRAM

C-14 depends on the untrusted creator of the video not knowing the sequence of motions, which we refer to as a *motion program*, before time t_b . There are two general ways to create the motion program: *unconstrained* and *constrained* programs. In an unconstrained program we use a hand-crafted sequence of motions over a wide-area including any type of motion a drone is capable of, such as yawing at different rates while changing the gimbal pitch. In contrast a constrained program uses constant motions, meaning that the translation and rotation does not vary during the motion. For instance if the drone translates forward, it does so without changing direction, or if it yaws, it does so at a constant rate. The advantage of a constrained program is that it is more amenable to sampling because the drone's rate of rotation or translation can be verified by looking only at randomly selected parts and assuming the motion is the same over that period. The advantage of an unconstrained program is that it may fit more naturally into systems where a particular path needs to be followed over a wide area, rather than a computed path.

4.1 Unconstrained Program

In unconstrained programs, we create motion programs that are sufficiently complicated such that guessing the program ahead of time would be improbable. As long as we only divulge the program to the video creator after t_b , we can show that the video was created after that time.

In Figure 2 we show such a motion program taken from a popular online drone flight planning system named Litchi [4]. Each plain numbered pin represents a *waypoint* for the drone to fly to, and each numbered pin containing a camera icon represents a *point of interest* for the drone to focus on during different parts of the mission. The point of interest has an altitude as well, creating yaw and gimbal pitch throughout the motion. The curves around waypoints represent the actual flight path to be taken to smooth out the drone's motion to and from the waypoint.

Such unconstrained motions are highly-complex, and thus there are a very large number of distinct possible drone paths. C-14 measures both the yaw between waypoints, as well as an average translation vector to verify the video. As we show in Section 7.6 even small deviations can be detected in the



Figure 2: A sample unconstrained mission from the Litchi flight planning software.

resulting video, ensuring that the number of distinct drone paths is very large.

A side benefit of using Litchi missions is that many users publicly post their flight plans and the resulting video from the drone, giving us a varied dataset to work with for our evaluation.

4.2 Constrained Program

The disadvantage of an unconstrained program is that the motion between waypoints is not a constant rate. So when verifying the unconstrained flights we have to sample more of the video, slowing down verification. An alternative to the unconstrained program is one that uses a small number of constrained, steady, and mathematically derived motions.

The constrained motion program is a sequence of motions, m_1, m_2, \dots deterministically derived from a number, N , with sufficient entropy where N is not known before t_b . To keep the system understandable m_n is always a combination of a rotation on one axis and/or a translation along one axis. An example motion would be flying in a circle while remaining pointed at a center point, which is a combination of a yaw rotation with translation to the side of the drone. This is commonly called a *hotpoint* motion: the drone circles, and points the camera at a point of interest. We use hotpoints as the key motion in the constrained program.

The hotpoint motion also must execute for a certain magnitude: the total number of degrees over the time period of the motion. To derive the number of degrees of rotation in the video, C-14 requires the *field of view* (FOV) of the camera. If the video creator supplies this parameter, it can scale the total amount of rotation perceived in the video. However, as we note in 6.3.3 we can use a typical FOV for a drone and achieve good results as drones typically use very similar cameras. One attack on the system is to statically crop or dynamically zoom the video. Zooming will create different views when looking at the object at different times (which happens many times in the constrained program)—this would be detectable. Cropping would scale all of the rotations, and finding a cropping value that makes all of the rotations match is difficult. We leave a thorough exploration of this issue as future work.

We also make the motions time-independent and only measure motions based on sequence of motions that occurred. This makes the system portable across drones with varying capabilities in speed and frame rate.

We have created a simple constrained program based on a single, high-entropy number, N , that would only be known after t_b . In this case we use the block hash from the Ethereum blockchain from a block that occurs shortly after t_b .

The drone flies a series of motions consisting of two parts: a motion to fly towards the point of interest to an inner radius and then back out again (called a fly in and out) followed by one hotspot motion for $angle$ degrees randomly chosen from 0 to 360 degrees along an outer radius, either clockwise or counter-clockwise (called a hotspot). The purpose of the translation is to provide a separation between hotspot motions. If there is no separation, the verifier cannot attribute the yaw to each required hotspot. During the fly in and out motion we only need to know that the drone did not appreciably yaw for some number of frames. The motion program ends with one more fly in and out to bookend the last hotspot.

All of the motions are done with the drone pointing its camera at one point of interest at the center of an circle, but with constant gimbal pitch. An example motion program is shown in Figure 3. The angle and direction of the hotspot motion is determined by the high-entropy number N by seeding a random number generator that produces a series of random numbers r_0, r_1, \dots . The motion program is a series of hotspot motions for a certain number of degrees, as: $hotpoint_i = (angle = r_{2i} \% 360, clockwise = abs(r_{2i+1}) \% 2 == 1 ? true : false)$.

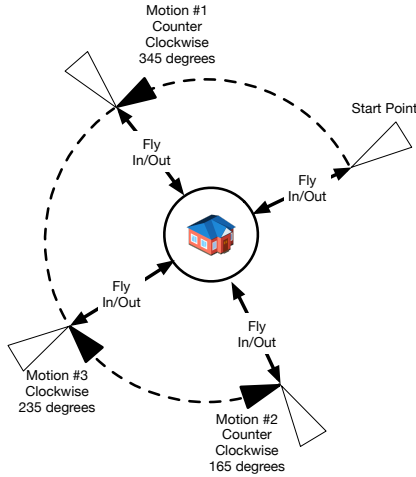


Figure 3: This figure shows the sequence of motions found in our example constrained motion program. The motions are a sequence of hotspot motions and fly in and out motions, all centered on a point of interest.

4.3 Probability of False Positives

An attacker could perform a brute-force attack: producing videos conforming to a random motion program and then

submitting them as authentic. While other stop-gaps would likely prevent a large number of videos being produced and submitted, it is valuable to know how likely it is that a random video would be accepted. Assuming that the number of motions N is known, a random attack can be done by picking N random rotations in $[-360, 360]$ degrees (-360 is a counter-clockwise full circle, and 360 is the same, but clockwise). As we show later, some tolerance (aka threshold) is needed to prevent large numbers of false-negatives. So given a true rotation R_t , a random rotation R_r is considered correct if it is in the range $[R_t - threshold, R_t + threshold]$. If we consider N rotations, the probability p_N for the N rotations to be correct is:

$$p_N = \left(\frac{2 * threshold}{720} \right)^N$$

In a constrained video with 4 rotations with a threshold of 20 degrees, that is approximately one chance in one hundred thousand.

We must also consider that the video creator can look back in time to find a motion program that matches a given video. For instance if we draw a random seed from a blockchain block, an attacker could generate as many motion programs as there are past blocks and claim that the video is from a matching block. One solution is to limit the granularity of time stamps to something larger than a block in the blockchain, such as the first block of the day. Then only a few thousand timestamps and blocks can be chosen from.

However, this motion program represents just the first possible constrained motion program. It is relatively easy to increase the entropy of the motion. For instance if the hotspot motion goes up or down in altitude and moves in or out (a spiral), we increase the possibilities by a factor of 4 and the probability of a false positive decreases to four in one hundred million.

In the unconstrained setting, we also need to consider the entropy of the translations. The probability for a random N -waypoint flight to match the flight plan becomes:

$$p_N = \left(\frac{2 * rotationThreshold}{720} \right)^{N-1} \times \left(\frac{2 * translationThreshold}{360} \right)^{N-1}$$

For a 10-waypoints flight the probability for a random video to match a given flight plan using a 30° rotation threshold and a 40° translation threshold is then smaller than 1 in 10^{17} .

5. VERIFIER

The goal of C-14 is to verify that the motion depicted in the video is consistent with a set of flight instructions given at the beginning of the flight. We leverage recent results from computer vision to estimate the camera motion, and thus the motion of the drone, based on the video. Using this estimated motion, we can verify that the drone has translated and rotated in the specified direction, and in the specified order.

The C-14 verifier takes in an untrusted video, a time-line description of when each motion element occurred (the metadata) and the timestamp claimed with the video. The verifier first ensures that the timestamp claimed for the video, t_b is consistent with the metadata. If that is true, then it must verify that the video matches the metadata. The verifier then produces a pass/fail determination based on the results of each test. We explain each step in detail here.

5.1 Verifying the Metadata

The metadata provided with the video describes the motions that should be contained in the video. The verifier first checks that these claimed motions are consistent with the timestamp claimed in the video. For the unconstrained video the process is straightforward, it must check that the flight plan in the metadata provided with the video is one that was not revealed until after time t_b . For the constrained videos, the verifier uses the timestamp to fetch the correct hash value from the blockchain. Given this hash value, the verifier computes the motion program using the same algorithm used by the video creator and ensures that it matches the motions contained in the metadata for the video.

Additionally the metadata describes where in the video each motion occurs. The verifier operates on the whole video with no gaps between motions. If gaps were allowed, then part of a motion would be ignored by the verifier, allowing the attacker to modify the motion without changing the video.

An alternative is to ignore the metadata entirely and use the computed motion in the video to recreate the metadata. For instance the verifier can look for where the drone stopped executing a yaw and started to go forward. This transition time would represent the change from a hotspot motion to a pure translation fly-in motion. However, this requires computing all, or a large part, of the motion estimation from the video, something that is computationally expensive. In Section 7.7 we show how sampling speeds up verifying the video. Note that using the metadata does not make verification weaker, it is just a labeling of where to find the data is—providing incorrect metadata will only make verification fail as part of one motion will be included in another.

5.2 Optical Flow and Motion Estimation

To compute the drone’s motion from the video we draw on recent results from computer vision in optical flow and motion estimation. Optical flow is the process of analyzing a video to show how pixels move from one frame to the next. If the scene is largely static (i.e., there are no moving objects) and the camera has no component of forward or backward motion, then the pixels move in a direction opposite that of the camera motion. For example if the camera moves to the right,

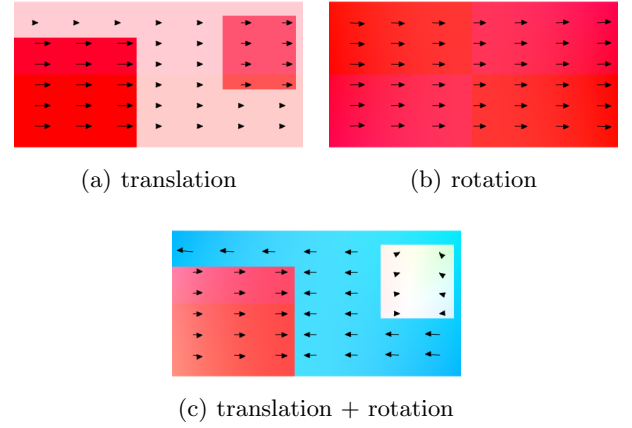


Figure 4: Optical flow due to camera motion of a static scene with static objects located at different depths. The hue of the background shows the angle of flow, and the intensity (or saturation) of the color shows the magnitude of motion.

then the pixels appear to move toward the left (see Figure 4(a)). The output of an optical flow algorithm is a two dimensional vector field that gives the magnitude and direction of the movement of each pixel in the scene.

Optical flow estimation and the closely related problem of camera motion estimation are both heavily-studied computer vision problems. New learning-based approaches have been able to improve speed by a significant margin [16,19,36,38,44] and have recently matched or even surpassed traditional methods in accuracy. The current best performing optical flow algorithm is PWC-Net [39], which uses a deep convolutional neural network to estimate the optical flow, which we use in our implementation.

Once optical flow has been estimated, we use it to estimate the motion of the camera. Camera motion occurs in three dimensions and can be broken into two components: rotation and translation. Consider the relationship between translation and rotation of the camera. When a camera rotates, it only changes what it is looking at, but objects at different depths do not appear to move in relation to each other. However, when translating, new parts of a scene become visible (disocclusions) or become hidden (occlusions), and objects that are closer appear to move faster than objects further away (see Figure 4). Results from photogrammetry and computer vision have shown that it is possible to disambiguate translation from rotation and thus discover how the pose of the camera is changing from frame to frame [12,13,32,46]. We use a recent camera pose estimator [13] to output a six-valued vector of the camera motion - the three rotation parameters pitch, yaw and roll and a three-dimensional unit vector defining the 3D translational motion direction. Note that this unit vector gives the direction, but not the magnitude (speed) of the translation direction.

The verifier computes the optical flow and motion estimation on the untrusted video which outputs an estimate of the camera motion on three translation axes and three rotation axes. It then matches the camera motion to the *ideal* motion the drone should have followed according to the flight plan or constrained program.

5.3 Translating Frame of Reference

Before making that match we must translate the video into the correct frame of reference. Recall that we describe everything from the *ideal drone* frame of reference described in Section 3. In the ideal drone frame of reference the only rotation is yaw (no pitch or roll), but the drone is free to translate on three axes. However, the video from the drone is taken from a frame of reference of the camera, which may be pitched by θ degrees (generally it is pointed down towards the ground) with respect to the ideal drone. In the constrained videos we fix θ at one value and that is provided with the metadata (or could be a globally fixed value). In the unconstrained videos θ needs to be generated as part of the flight plan. This angle is typically not constant. The camera does not roll with respect to the horizon as it has a gimbal, and thus the camera and the ideal drone have no roll. The camera is also set to *yaw follow* mode which means that the yaw generally matches the ideal drone with some elasticity. In the interests of space, we omit the details of this translation here and will include it in the documentation of the source code at publication.

5.4 Verifying Motion in Unconstrained Videos

Figure 5 shows an example of the output of the camera motion in the drone frame of reference compared to the drone motion from the flight plan.

To verify unconstrained videos, we check that between each pair of waypoints the camera motion output and the flight plan match. Specifically we check that (i) the drone’s total yaw and (ii) the average error of the angle between the two translation vectors, are both within some threshold.

Let’s define two sequences of motions expressed in the ideal drone frame of reference: (i) the sequence of rotations and translations at each frame of video from the camera motion estimator and (ii) the same sequence of rotations and translations from the flight plan. While the flight plan does not contain speed information, we can still make a correspondence between the output of the camera motion and the flight plan by matching the waypoints in each and then interpolating between points in the flight plan. This requires knowing where in the video the drone reached each waypoint, which we discuss in Section 6.1.

The verifier then sums the yaw in the flight plan between each pair of waypoints and compares it to the sum of the yaw between two waypoints from camera

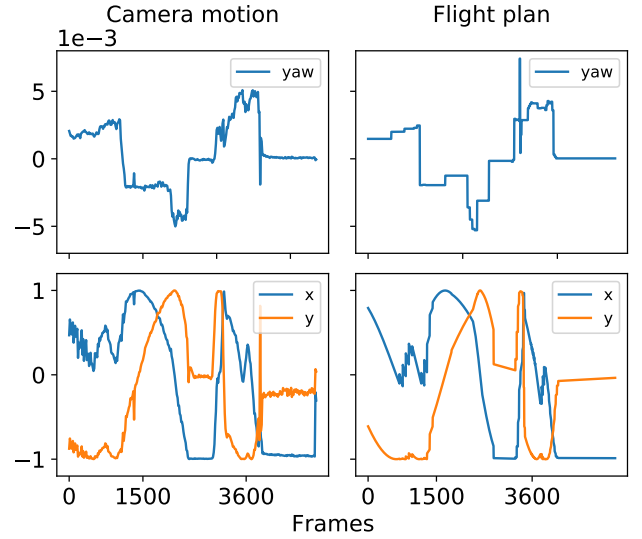


Figure 5: Comparison of the camera motion (left) and the Litchi flight plan (right). The camera motion has been computed from a 3-minute video at a resolution of 384*216 and a frame rate of 2 FPS.

motion. If the difference of the two sums are within a threshold for all waypoints, then the verifier is satisfied.

To verify the translation, the verifier averages the translation vectors in the flight plan between two waypoints and compares it to the averages of the translation vectors between two waypoints from the camera motion. The verifier then checks if the angle between the two resulting vectors is smaller than a threshold. Note that the average of the translations does not have any physical meaning because each translation vector is expressed in the ideal drone frame of reference at their respective times.

5.5 Verifying Motion in Constrained Videos

The motion program instructs the drone to complete a total amount of yaw during the hotpoint. To find the yaw, the verifier samples N frames of motion, chosen uniformly randomly in the interval, fits a curve to the samples using piecewise linear fit [5], and then integrates the curve. It then checks if the total yaw is within a threshold of the target and a large percentage of the yaw samples are the correct sign (for a counterclockwise motion the yaw should be negative). We discuss all of the thresholds we used in Section 6.

During this motion the drone is also flying in a positive or negative y -axis direction, which corresponds to a counter-clockwise or clockwise hotpoint respectively. Recall this is in the drone frame of reference, so while doing a hotpoint the drone yaws between frames and moves on the y -axis only. Similar to yaw, a large percentage of the samples must be the correct sign.

Recall that the fly in and out motion is only there to bookend the hotspot motions. To verify the fly in/out motion, the verifier checks that there is no appreciable yaw for a certain number of frames and then the verifier accepts this as correct.

5.6 Sampling

To verify a 133 second video (3993 frames) at 4K UHD resolution (3840x2160) and 29.97 fps, the computation time required for optical flow and camera motion is 26 hours on a high-end Xeon processor and a modern NVIDIA GPU.

We use four techniques to reduce the amount of computation needed to approximately 1 minute (a factor of more than 1500): video compression, frame skipping, spatial sampling, and temporal sampling. Video frame compression simply reduces the resolution of the frames. Verifying lower resolution videos is faster, but may prevent optical flow from recognizing corresponding pixels across subsequent frames. Frame skipping reduces the frames per second of the video, reducing computation as well, but skipping too many frames poses similar difficulties for optical flow. However, we have found skipping some frames generally has a *positive* effect as it has a smoothing effect on the estimated motion (up to a point where optical flow breaks down). Further, the camera motion algorithm can sample the optical flow by taking the center point from a square grid spaced at N pixels. Larger grid squares reduce computation, but only up to a point.

For temporal sampling we only examine the motion in discrete parts of the video frames that supposedly contain that motion. This is similar to skipping frames, and the verifier samples after reducing the frame rate, but it does so using uniform random sampling to prevent an attacker from exploiting the reduced sampling.

We show in the results that for both unconstrained and constrained videos we can reduce the resolution by a factor of 144, skip up to 15 frames, and sample the optical flow in a 7x7 pixel grid. For constrained videos we can process just 4% of the frames (skip 15 frames and sample 60% of those) over the whole video. This combination of techniques reduces the computation time to less than three minutes.

6. IMPLEMENTATION

Our implementation of C-14 consists of two parts: a method for capturing videos and metadata from flights, and a system for verifying that videos match the intended motion program. We gathered videos from two sources: (i) public videos of drone flights with flight plans and (ii) an implementation of a motion program using the DJI SDK [2]. At publication we will make all of the source code of the drone program and verifier public, as well as links to all of the data sets we use in the evaluation.

6.1 Unconstrained Videos

To collect unconstrained motion programs and videos we take advantage of a popular drone flight planning application (Litchi) which allows users to publicly post their flight plans on the Litchi website and the resulting video on YouTube. Using videos from third parties helps eliminate potential bias in data collection and gives us access to a large variety of scenes (rural, nature, cities, etc.), lighting conditions, and DJI drone models.

A disadvantage of this data set is that we do not have the metadata (GPS, heading, speed, etc.) so we do not know when in the video the drone executed each movement (such as flying around a waypoint). In a deployed C-14 system we would have access to the metadata as it would be submitted to the verifier with the video. However, by using the camera motion results obtained from a full video, we can manually label when each motion starts.

The Litchi flight plan only contains the waypoints, heading, and gimbal position. To translate this into a motion in the ideal drone frame of reference, we first obtain the drone poses in the world frame of reference using Virtual Litchi Mission [7]. The drone poses are a list of timestamps describing the longitude, latitude, altitude, heading, tilt and roll of the camera. From those poses, the motions in the ideal drone frame of reference can be derived. By definition, there is no tilt and no roll in the ideal drone frame of reference. The yaw corresponds to the difference of heading between two consecutive timestamps.

$$yaw_{(t,t+1)} = heading_{(t+1)} - heading_{(t)}$$

The translation along the z-axis is the difference of the altitudes.

$$t_{(y_{t,t+1})} = altitude_{(t+1)} - altitude_{(t)}$$

To compute the translations along the x-axis and the y-axis we use the geodesic distance with the WGS-84 ellipsoid model. The translation along the x-axis is given by the distance between the point $(lat_t, long_{t+1})$ and the point $(lat_t, long_t)$. The translation along the y-axis is given by the distance between the point $(lat_{t+1}, long_t)$ and the point $(lat_t, long_t)$. To express the translation vector in the ideal drone frame of reference, we rotate the translation vector by the heading of the drone. We then normalize the translation vectors to match the normalized translation vectors coming from the camera motion estimator.

6.2 Constrained Videos

To collect constrained videos we implemented the C-14 motion program using the DJI Mobile SDK, which is compatible with many DJI drones. We implemented the motion program using DJI's mission control API that uses a mobile device to program the drone with a series of mission elements (waypoint, hotspot, etc.).

One disadvantage of using the mission control API is that transitions between motions (such as a hotpoint to a fly in and out motion) have a multi-second pause. Future systems could eliminate this pause by more directly controlling the drone through the virtual stick [3].

Our interface to the DJI SDK was built into a React Native application using a wrapper for the DJI SDK [6]. We contributed a number of changes to the library to implement capabilities required for C-14. Our C-14 program runs on Android and connects to the drone via WiFi and the DJI remote controller. The app, shown in Figure 6 has functions to create new motion programs for a particular location chosen on a map and run the program on the drone.



Figure 6: This figure shows a screen from the constrained video application running on an Android device. We used this app to collect videos for the evaluation and to provide a starting point for future constrained programs.

While working on the system we discovered two issues: the hotpoint missions in the DJI SDK do not accurately fly the given number of degrees, instead flying too many or too few degrees by many 10s of degrees. As C-14 depends on the flight to be accurate, we re-implemented the hotpoint mission to measure the GPS location of the drone and stop just short of the required number of degrees. Future systems could make this even more accurate by computing the optical flow and camera motion directly on the drone to ensure accurate compliance with the required motion program.

6.3 Verifier

The verifier consists of several stages: sampling and compression, optical flow computation, motion estimation, and verification. We implemented the verifier as a Python script and it invokes an implementation of the optical flow estimator, PWC-Net, that has been opti-

mized for GPUs [29]. We use a Matlab implementation of the camera motion estimator [13] obtained from the authors.

6.3.1 Sampling and Compression

As described in Section 5.6 we compress the videos and reduce the frame rate to speed up verification. For compression we use the `resize()` function of OpenCV with the `INTER_AREA` interpolation algorithm. To reduce the frame rate, we maintain one frame among every s frames. For example, if the skip rate is 5, we only maintain the frame whose indexes are the multiples of 5, i.e. frame #0, #5, #10.

6.3.2 Segmentation

The verifier looks at frames and samples from the entire submitted video to ensure that it conforms to the target program. The metadata describes when each motion starts and ends (fly in, fly out, hotpoint). We do not allow gaps between the motions as it would allow an attacker to possibly modify the results. For instance, an attacker could say that a hotpoint or turn started later than it does in the video which could reduce the yaw value to match a target. Using the metadata we separate the video into segments: for unconstrained videos this is the time between waypoints, and for constrained videos it is each fly in, fly out, and hotpoint motion.

6.3.3 Verification

After segmentation, the verifier checks if each segment conforms to its corresponding motion as described in Section 5. For constrained videos we only sample parts of the motions as described in Section 5.6. The evaluation section shows how many samples are needed to achieve good results. The verifier computes the optical flow for all of the samples first and then computes the camera estimation. If any of the segments fail the verifier fails the video.

As discussed in Section 4.2 the verifier requires knowing the FOV of the drone camera—it is specified as a ratio of the focal length to the sensor width. For all of the videos we use a parameter of 0.82, which we measured from the Mavic Air. However, we did not change this parameter for the unconstrained videos because we do not know the model of the drone used. Fortunately, the system is relatively insensitive to the FOV parameter because most drones have a similar camera.

Also, we must add some amount of tolerance to the verifier to account for inaccuracies in the camera motion estimation process and sampling noise. The angle tolerance is measured in the evaluation section. During a hotpoint we check that 60% of the samples have the correct Y -axis sign. For the fly in and out, we check that the integral of the yaw is no greater than 15 degrees.

7. EVALUATION

To evaluate C-14 we collected two data sets: unconstrained videos taken from Litchi and constrained videos taken using our custom drone program. We downloaded 10 unconstrained flight plans and videos from Litchi and YouTube respectively. The videos with links to the flight plans are shown in Figure 7 and were collected internationally with a variety of drones. We also collected eight unconstrained videos ourselves. For constrained videos we collected videos from a number of different settings, at different altitudes, and different radius and number of motions, shown in Figure 8. We collected the constrained videos using a DJI Mavic Air drone. This is a relatively inexpensive drone (\$900USD) with a three-axis gimbal, a 4K UHD (3840x2160) 30 fps camera, and capable of speeds of 30 km/h in obstacle avoidance mode.

Name	FlightID	Length/Motions
Church	j6uTc0Qvha	3:27/13
Vegetation	mQUaT4UHLA	3:01/16
Small town	bNqRdjSFmo	3:44/10
Garden	qBV1h0veQ2	3:46/15
Village	bHg7fNYTSW	3:27/13
Ruins	u6UgiEDrp5	4:54/10
School	cSBHZth7L2	2:22/7
Ice hockey	mVCQVhgy0c	1:47/6
Tree	IFaQrqidsy	2:08/8
Ice rink	k7e4Hmi9Gm	1:06/7
House(x8)	withheld	approx. 0:56/8

Figure 7: This table lists the 18 flights used in our unconstrained data set. We did not create these flights and videos (except the House video). The flight plans can be obtained from <http://flylitchi.com/hub?m=FlightID>. The videos can be obtained from the Litchi links by double clicking on the yellow "eyeball" icon. We will provide a full data set upon publication.

Name	Altitude(m)	Radius	Length(s)/Motions
Forest Road	55	15/30	88/2
Creek	55	15/30	86/2
House A	55	15/30	119/2
Roadway(x2)	40	15/30	94,102/2
House B(x2)	80	20/40	91,106/2
Barn	40	15/30	78/2
Soccer Field	30	15/30	95/2
Farm Field	30	15/30	100/2
Snow House(x2)	30	15/30	101,106/2
Library(x3)	40	15/30	93,190,184/4
School(x4)	40	15/30	171,196,189,173/4
River(x3)	55	15/30	176,156,157/4
RecCenter(x2)	30	15/30	177,176/4
Parking lot(x2)	30	15/30	216,212/4

Figure 8: This table lists the 26 flights used in our constrained data set with their altitude, the radius of the inner and outer circles, the length of the video in seconds, and number of individual motions(fly in/out+hotpoint). Roadway and House B were at the same location with the same program. The remaining multiple videos used different programs.

We ran the experiments on a GPU cluster with various CPUs and GPUs. All of the timing experiments were completed on a machine with two Xeon E5-2620 v3 2.40 GHz CPUs, 256G of RAM, and an NVIDIA TITAN X GPU.

7.1 Resizing and Spatial Sampling

All of the sampling methods create a trade-off between (i) the accuracy of the camera motion estimator and (ii) the run time of the system. To evaluate the accuracy of the camera motion estimator, we use a constrained video and it's associated metadata, which records the drone's location via GPS, as well as it's heading via an onboard compass, and speed. We compute the Mean Squared Error (MSE) of the drone's yaw at each frame between the camera motion estimation and the metadata. The results are shown in Figure 9.

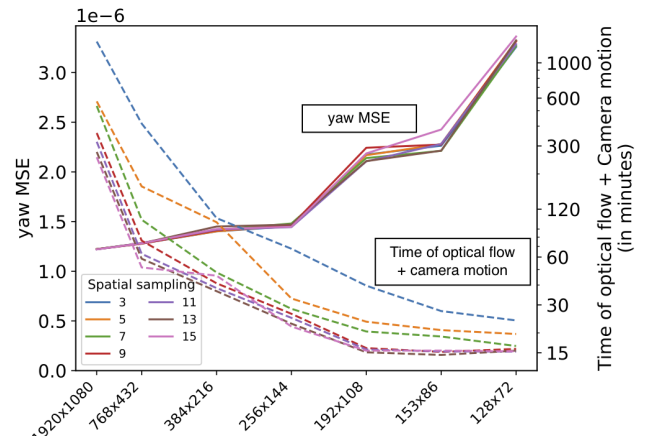


Figure 9: This shows the comparison of the computational time and Mean Squared Error of the yaw between the metadata and the camera motion estimator for different video resolutions and different spatial sampling. Each data point is the average of three runs on a 2-minute video.

The results show that the computational time increases exponentially with the resolution of the video, while the error decreases exponentially. Using these results as a guide, we choose to resize the videos to 320x180 (a factor of 144 fewer pixels) for the rest of the results in the paper.

Similarly, we evaluate the impact of the camera motion spatial sampling. As shown in Figure 9, all the spatial sampling values seem to give very similar accuracy results. Based on measuring the impact on the runtime of camera motion (not shown) we choose a spatial sampling of 7 as beyond that the motion estimator does not run appreciably faster.

7.2 Frame Rate

We collected videos at 30 fps, however accurate optical flow does not require all of the frames and fewer frames

means less processing time. Using the same video as before we compute the MSE between the camera motion estimator and the metadata. The results are shown in Figure 10.

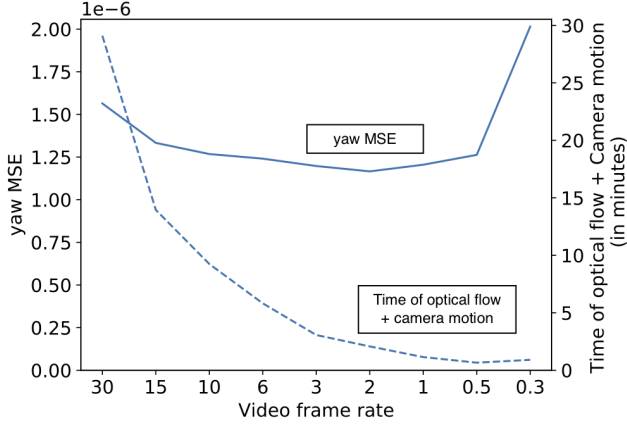


Figure 10: This shows the comparison of the yaw Mean Squared Error and the computational time for different frame rates. We used 256x144 video and a spatial sampling of 7. Each data point is the average of three runs on a 2-minute video.

In addition to computational time benefits, skipping frames reduces the MSE up to a point. This is because skipping frames smooths the motion between sample points by measuring a larger motion. However, once the system skips too many frames optical flow has a harder time making pixel correspondences between frames and the MSE increases. Based on these results we choose to use a frame rate of 2 fps.

7.3 Sampling Rate

In constrained videos we additionally use sampling to reduce the amount of computation required. To find a reasonable rate we sample a portion of the frames (after reducing the frame rate), and compute the error between the yaw estimation from the camera motion estimator and the motion program target angle. A histogram of the results is shown in Figure 11. The results show that below 60% sampling the yaw error grows outside of the bounds of $[-20, 20]$. Thus we fix the sampling rate at 60%.

7.4 Yaw Error

A key component of verifying videos is checking the total yaw of the drone during a hotspot motion in constrained videos, or between waypoints in unconstrained ones. We use all of the hotspot motions (except House) from the dataset in Figure 7 (104 motions) and Figure 8 (38 motions). We use the metadata as ground truth for the constrained videos and the flight plan for the unconstrained videos. The distribution of error is shown in Figure 12.

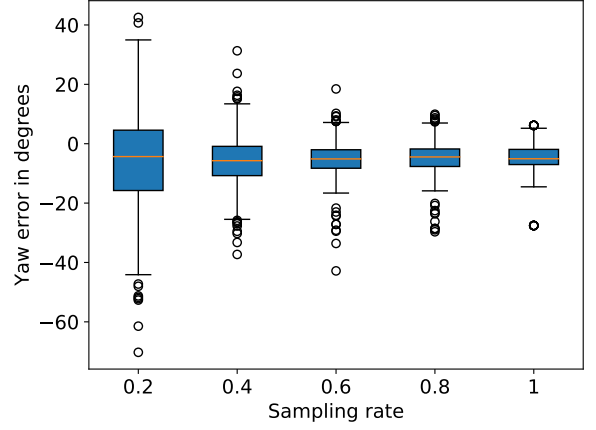


Figure 11: This shows the yaw error for various sampling rates. For a sampling rate larger than 60%, the yaw errors remain in the same range of values. A smaller sampling rate makes the yaw error increase significantly.

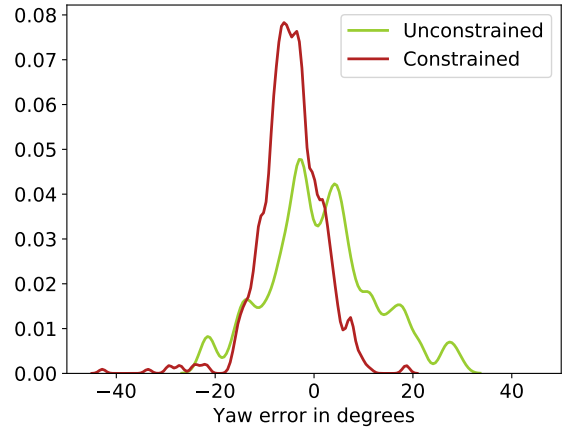


Figure 12: This shows the distribution of yaw errors for the constrained and unconstrained videos.

The results confirm that thresholds on the order of $[-20, 20]$ for constrained videos and $[-30, 30]$ for unconstrained will ensure that a high percentage of yaw angles will be classified as correct. Since the unconstrained videos allow more complicated camera motions, and the ground truth for the yaw is the flight plan, rather than the metadata, the camera motion estimation generates somewhat larger errors.

Looking further into the data we find that all of the errors in Figure 11 and 12 that fall outside of a threshold of 20 degrees are due to a single hotspot in five different videos: Forest Road, two River videos, one Rec-Center, and one School video. We believe that optical flow has trouble making accurate pixel correspondences due to the textures in those videos (leafless trees in the first three and lots of snow in the fourth and fifth).

However, if we change the frame rate of these videos to 3 frames per second and 80% sampling, the error decreases dramatically (and well within 20 degrees). In a deployed system we can rerun any detected negatives with higher sampling rates, frame rates, and lower compression to truly verify negatives and eliminate such a case, and we apply this technique in the next section.

7.5 False Negative Rate

We evaluate the false negative rate with different thresholds in both the constrained and unconstrained settings. A false negative occurs when a video is created using a motion program or flight plan and it is rejected by the verifier. For constrained videos we sample at various rates, processing each video 10 times as the sampling is random. We plot the resulting false negative rate in Figure 13 for various yaw thresholds and sampling rates. The results show that for constrained videos we can achieve a 0 false negative rate at a sampling rate of 60% a frame rate of 2fps (with the exception of the five videos mentioned previously at 80% sampling and 3fps), and a threshold of 20 degrees for yaw.

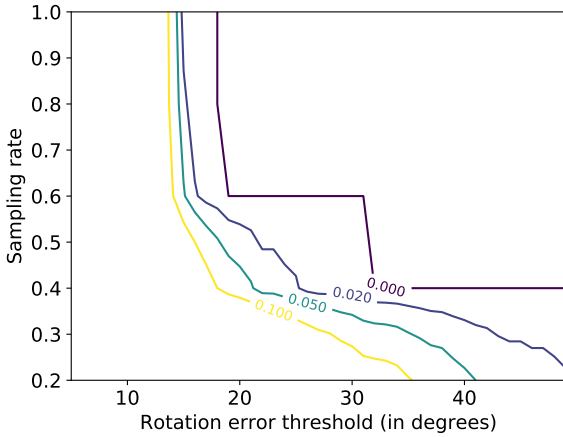


Figure 13: This shows the false negative rate with respect to the rotation error threshold and the sampling rate.

Verifying unconstrained videos does not use sampling and is therefore deterministic. We plot the false negative rate in Figure 14 using various thresholds for both yaw and translation vectors. With a threshold of 30 degrees for the rotations and a threshold of 35 degrees for the translations, the system achieves a 0 false negative rate.

7.6 False Positive Rate

For constrained videos we confirm that given the motion programs in the dataset in Figure 8, each only matches the video it was used for. However, given the size of the dataset, the analytical evaluation of the false positive rate for constrained videos is a better measure (see Section 4.3).

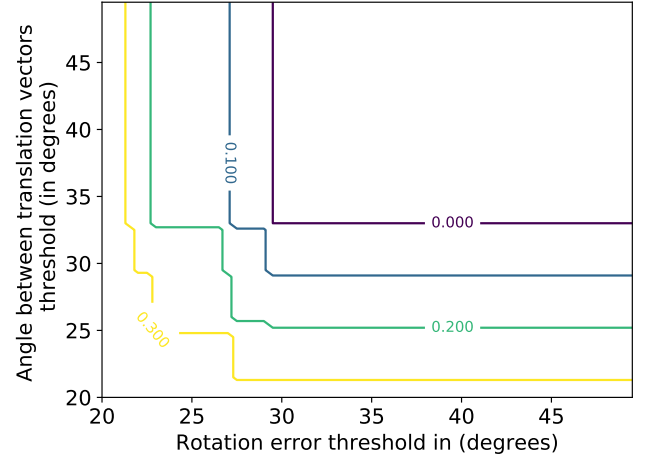


Figure 14: This figure shows the false negative rate at a particular rotation threshold and translation threshold.

For unconstrained videos a false positive will only occur if a video from some flight plan A matches some other random flight plan B. We evaluate the chance of this happening in Section 4.3 as well. However, the creator of the video might not want to create a completely random flight plan, but rather modify an existing one by “enough” such that previous videos will not match the new plan. We evaluate what “enough” is by modifying one waypoint in the flight plan. We use the chosen thresholds for unconstrained video i.e. 30 degrees for the rotation angle and 35 degrees for the translation angle. In Figure 15 we show the false positive data points in terms of how large the maximum rotation difference between the computed angle from the camera motion and the target angle from the flight plan are. This demonstrates that modifying a single waypoint in the flight plan to change the yaw of the drone by 30 degrees or more will eliminate false positives.

We also confirm that subtle changes in an unconstrained flight plan creates a distinguishable video. We programmed our drone with four Litchi flight plans, approximately 340 meters of flight, and took two videos from each plan. The second flight plan differed from the first by moving just one of the waypoints by 17.5 meters. The third moved the same waypoint by 56 meters and the fourth by 63.5 meters. The second flight plan produces a video that is a false positive with the first flight plan, but the other two flight plans verify as true negatives. This shows that moving just waypoint by 50 meters was sufficient to prevent false positives.

7.7 Computational Time

We evaluate the computational time using two Xeon CPUs and a TITAN X GPU. Processing a 59 second unconstrained video takes 91 seconds and a 190 second constrained video (with sampling 60%) takes 158

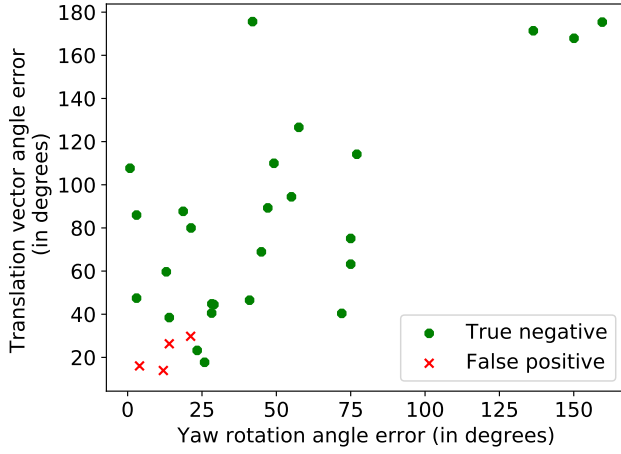


Figure 15: This figure demonstrates how much a flight plan must change to eliminate false positives. If a single waypoint is changed to produce 30 degrees more or less yaw and 35 degrees of average translation vector direction, no false positive will occur.

seconds. 60% of the time is spent on camera motion estimation, 35% on the optical flow and 5% on the compression. Verifying the motion is within the tolerance takes negligible time.

7.8 Overhead

One concern is how long it takes the drone to complete the motion program. An assumption with unconstrained videos is that the motion program is part of the flight plan that would have been used anyway, so there is no overhead for incorporating C-14.

For the constrained videos there is also an assumption that the motions sequence’s center point is a subject of interest, and therefor the motions are not purely overhead. However if we conservatively take the entire motion sequence as overhead the overhead is $N \cdot (180/R + 2 \cdot D/V)$, where there are N motions, the average rotation is 180 degrees, the copter completes the hotspot at R degrees/sec and a fly in/out motion covers D meters between the outer and inner radius at V m/s. Our copter can hotspot at approximately $R = 10 \text{ degrees/sec}$ at a radius of 30m and can fly at $V = 8.3 \text{ meters/s}$. So a flight with 4 motions at an outer and inner radius of 30m and 15m, will take $4 \cdot (180/10 + 2 \cdot 15/8.3) = 86$ seconds. However, the drone must accelerate to full speed, does not stop on a dime, and due to limitations of the DJI SDK pauses for many seconds between motions. Thus a 4 motion sequence takes approximately 3 minutes (see Figure 8), which is 15% of the Mavic Air’s flight time.

8. RELATED WORK

Drones have inspired a great deal of recent work in the mobile systems community, including testbeds [8], control algorithms [15,24], and detecting the presence of

a drone [28]. However, we are unaware of any system that attempts to place a timestamp on a video using techniques similar to those in C-14. Perhaps the closest example might be so-called “proof of life” videos where a subject is holding a recent newspaper in a video. In contrast C-14 works from drones over wide areas with no special augmentation of the environment.

A similar, but different property is that of “liveness”: the provider of a video is the original creator of a video. In two systems, Vamos [33] and Movee [34], Rahman et. al. seek to verify a user’s claim that they created a particular video. When creating the video a user also records readings from the inertial sensors, such as the accelerometer. When providing proof to a trusted third party, the user provides that sensor stream, which can be matched to the movements in the video. In contrast we are not proving ownership, we show that a video was taken in a particular time window.

An alternative to using optical flow is to use monocular Simultaneous Localization and Mapping (SLAM) [18,26]. SLAM algorithms estimate the drone’s position while simultaneously estimating a map of the unknown 3D environment - while none of these are known beforehand. This is also often referred to as the chicken-and-egg-problem. Instead, we solely estimate the drone’s camera rotation and translational motion direction to verify unmodified drone videos - an easier problem.

Providing assurance for videos does depend on the videos being unaltered, specifically not being spliced together from multiple videos. The more dynamic the pattern, the more and more difficult it becomes to gather enough videos to create a convincing fake video. While much has been made recently of “deep fake” videos, these have largely been applied to creating videos that change what a person is saying, such as in Face2Face [40]. Given the high value of such fakes, techniques have emerged that can detect such alterations [25,27,42,45]. Work has also been conducted on detecting computer graphics [35]. We are not aware of such systems applied to drone videos, but once fakes emerge, we believe researchers will combat them with similar techniques.

9. CONCLUSIONS

We consider C-14 to be an important first step in opening up the problem to defenses, forensics, and anti-forensics research [14]. While it significantly raises the bar for assuring the age of drone videos, the basic technique can be used for other purposes, such as ensuring the quality of videos resulting from drone flights, or reconstructing a flight plan from videos with no other information. We hope that others will benefit from the dataset and results.

10. REFERENCES

- [1] Using Drones to Shoot War Zones. <https://petapixel.com/2018/02/20/using-drones-shoot-war-zones/>, 2018.
- [2] DJI Mobile SDK. <https://developer.dji.com/mobile-sdk/documentation/introduction/flightController-concepts.html>, 2019.
- [3] DJI Mobile SDK. <https://developer.dji.com/mobile-sdk/documentation/introduction/component-guide-flightController.html#virtual-sticks>, 2019.
- [4] Litchi. <https://flylitchi.com/>, 2019.
- [5] Python pwlf. <https://pypi.org/project/pwlf/>, 2019.
- [6] React Native Wrapper Library For DJI Mobile SDK. <https://github.com/Aerobotics/react-native-dji-mobile>, 2019.
- [7] Virtual Litchi Mission. <https://mavicpilots.com/threads/virtual-litchi-mission.31109/>, 2019.
- [8] Mikhail Afanasov, Alessandro Djordjevic, Feng Lui, and Luca Mottola. Flyzone: A testbed for experimenting with aerial drone applications. In *Proceedings of the 17th Annual International Conference on Mobile Systems, Applications, and Services*, pages 67–78. ACM, 2019.
- [9] Javad Abbasi Aghamaleki and Alireza Behrad. Inter-frame video forgery detection and localization using intrinsic effects of double compression on quantization errors of video coding. *Signal Processing: Image Communication*, 47: 289–302, 2016.
- [10] Jamimamul Bakas and Ruchira Naskar. A digital forensic technique for inter-frame video forgery detection based on 3d cnn. In *International Conference on Information Systems Security*, pages 304–317. Springer, 2018.
- [11] Jamimamul Bakas, Ruchira Naskar, and Rahul Dixit. Detection and localization of inter-frame video forgeries based on inconsistency in correlation distribution between haralick coded frames. *Multimedia Tools and Applications*, 78(4): 4905–4935, 2019.
- [12] Pia Bideau and Erik Learned-Miller. It’s moving! a probabilistic model for causal motion segmentation in moving camera videos. In *European Conference on Computer Vision (ECCV)*, 2016.
- [13] Pia Bideau, Aruni RoyChowdhury, Rakesh R Menon, and Erik Learned-Miller. The best of both worlds: Combining cnns and geometric constraints for hierarchical motion segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 508–517, 2018.
- [14] Rainer Böhme and Matthias Kirchner. Counter-forensics: Attacking image forensics. In *Digital Image Forensics*, pages 327–366. Springer, 2013.
- [15] Endri Bregu, Nicola Casamassima, Daniel Cantoni, Luca Mottola, and Kamin Whitehouse. Reactive control of autonomous drones. In *Proceedings of the 14th Annual International Conference on Mobile Systems, Applications, and Services*, pages 207–219. ACM, 2016.
- [16] Thomas Brox and Jitendra Malik. Large displacement optical flow: descriptor matching in variational motion estimation. *IEEE transactions on pattern analysis and machine intelligence*, 33(3): 500–513, 2011.
- [17] Sintayehu Dehnie, Taha Sencar, and Nasir Memon. Digital image forensics for identifying computer generated and digital camera images. In *Image Processing, 2006 IEEE International Conference on*, pages 2313–2316. IEEE, 2006.
- [18] Jakob Engel, Thomas Schöps, and Daniel Cremers. Lsd-slam: Large-scale direct monocular slam. In *European conference on computer vision*, pages 834–849. Springer, 2014.
- [19] Junhwa Hur and Stefan Roth. Joint optical flow and temporally consistent semantic segmentation. In *European Conference on Computer Vision*, pages 163–177. Springer, 2016.
- [20] Suyog Dutt Jain, Bo Xiong, and Kristen Grauman. Fusionseg: Learning to combine motion and appearance for fully automatic segmentation of generic objects in videos. In *2017 IEEE conference on computer vision and pattern recognition (CVPR)*, pages 2117–2126. IEEE, 2017.
- [21] Shan Jia, Zhengquan Xu, Hao Wang, Chunhui Feng, and Tao Wang. Coarse-to-fine copy-move forgery detection for video forensics. *IEEE Access*, 6:25323–25335, 2018.
- [22] Staffy Kingra, Naveen Aggarwal, and Raahat Devender Singh. Inter-frame forgery detection in h. 264 videos using motion and brightness gradients. *Multimedia Tools and Applications*, 76(24):25767–25786, 2017.
- [23] Vincent Lenders, Emmanouil Koukoumidis, Pei Zhang, and Margaret Martonosi. Location-based trust for mobile user-generated content: applications, challenges and implementations. In *Proceedings of the 9th workshop on Mobile computing systems and applications*, pages 60–64. ACM, 2008.
- [24] Wenguang Mao, Zaiwei Zhang, Lili Qiu, Jian He, Yuchen Cui, and Sangki Yun. Indoor follow me drone. In *Proceedings of the 15th Annual International Conference on Mobile Systems, Applications, and Services*, pages 345–358. ACM, 2017.
- [25] Falko Matern, Christian Riess, and Marc Stamminger. Exploiting visual artifacts to expose

- deepfakes and face manipulations. In *2019 IEEE Winter Applications of Computer Vision Workshops (WACVW)*, pages 83–92. IEEE, 2019.
- [26] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *IEEE transactions on robotics*, 31(5):1147–1163, 2015.
- [27] Huy H Nguyen, Junichi Yamagishi, and Isao Echizen. Capsule-forensics: Using capsule networks to detect forged images and videos. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2307–2311. IEEE, 2019.
- [28] Phuc Nguyen, Hoang Truong, Mahesh Ravindranathan, Anh Nguyen, Richard Han, and Tam Vu. Matthan: Drone presence detection by identifying physical signatures in the drone’s rf communication. In *Proceedings of the 15th Annual International Conference on Mobile Systems, Applications, and Services*, pages 211–224. ACM, 2017.
- [29] Simon Niklaus. A reimplement of PWC-Net using PyTorch. <https://github.com/sniklaus/pytorch-pwc>, 2018.
- [30] Feng Pan, JiongBin Chen, and JiWu Huang. Discriminating between photorealistic computer graphics and natural images using fractal geometry. *Science in China Series F: Information Sciences*, 52(2):329–337, 2009.
- [31] Fei Peng, Jiao-ting Li, and Min Long. Identification of natural images and computer-generated graphics based on statistical and textual features. *Journal of forensic sciences*, 60(2):435–443, 2015.
- [32] Clement Pinard, Laure Chevalley, Antoine Manzanera, and David Filliat. Learning structure-from-motion from motion. In *The European Conference on Computer Vision (ECCV) Workshops*, September 2018.
- [33] Mahmudur Rahman, Mozghan Azimpourkivi, Umut Topkara, and Bogdan Carbunar. Video liveness for citizen journalism: Attacks and defenses. *IEEE Transactions on Mobile Computing*, 16(11):3250–3263, 2017.
- [34] Mahmudur Rahman, Umut Topkara, and Bogdan Carbunar. Seeing is not believing: Visual verifications through liveness analysis using mobile devices. In *Proceedings of the 29th Annual Computer Security Applications Conference*, pages 239–248. ACM, 2013.
- [35] Nicolas Rahmouni, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. Distinguishing computer graphics from natural images using convolution neural networks. In *2017 IEEE Workshop on Information Forensics and Security (WIFS)*, pages 1–6. IEEE, 2017.
- [36] Jerome Revaud, Philippe Weinzaepfel, Zaid Harchaoui, and Cordelia Schmid. Epicflow: Edge-preserving interpolation of correspondences for optical flow. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1164–1172, 2015.
- [37] Stefan Saroiu and Alec Wolman. Enabling new mobile applications with location proofs. In *Proceedings of the 10th workshop on Mobile Computing Systems and Applications*, page 3. ACM, 2009.
- [38] Laura Sevilla-Lara, Deqing Sun, Varun Jampani, and Michael J Black. Optical flow with semantic segmentation and localized layers. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3889–3898, 2016.
- [39] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. In *CVPR*, 2018.
- [40] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2387–2395, 2016.
- [41] P. Tokmakov, K. Alahari, and C. Schmid. Learning motion patterns in videos. In *CVPR*, 2017.
- [42] Weihong Wang and Hany Farid. Exposing digital forgeries in video by detecting double quantization. In *Proceedings of the 11th ACM workshop on Multimedia and security*, pages 39–48. ACM, 2009.
- [43] John Wihbey. The drone revolution - uav-generated geodata drives policy innovation. *Land Lines Magazine*, pages 14–21, October 2017.
- [44] Jonas Wulff, Laura Sevilla-Lara, and Michael J. Black. Optical flow in mostly rigid scenes. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [45] Xin Yang, Yuezun Li, and Siwei Lyu. Exposing deep fakes using inconsistent head poses. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8261–8265. IEEE, 2019.
- [46] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1851–1858, 2017.
- [47] Zhichao Zhu and Guohong Cao. Applaus: A privacy-preserving location proof updating system for location-based services. In *2011 Proceedings IEEE INFOCOM*, pages 1889–1897. IEEE, 2011.