# Vision and Learning for Deliberative Monocular Cluttered Flight

Debadeepta Dey[1], Kumar Shaurya Shankar[2], Sam Zeng[3], Rupesh Mehta[4], M. Talha Agcayazi[5],
Christopher Eriksen[6], Shreyansh Daftry[7], Martial Hebert[8], J. Andrew Bagnell[9]

*Abstract*—**Cameras provide a rich source of information while being passive, cheap and lightweight for small and medium Unmanned Aerial Vehicles (UAVs). In this work we present the first implementation of receding horizon control, which is widely used in ground vehicles, with monocular vision as the only sensing mode for autonomous UAV flight in dense clutter. We make it feasible on UAVs via a number of contributions: novel coupling of perception and control via relevant and diverse, multiple interpretations of the scene around the robot, leveraging recent advances in machine learning to showcase anytime budgeted cost-sensitive feature selection, and fast non-linear regression for monocular depth prediction. We empirically demonstrate the efficacy of our novel pipeline via real world experiments of more than 2 kms through dense trees with a quadrotor built from off-the-shelf parts. Moreover our pipeline is designed to combine information from other modalities like stereo and lidar as well if available.**

## I. INTRODUCTION

Unmanned Aerial Vehicles (UAVs) have recently received a lot of attention by the robotics community. While autonomous flight with active sensors like lidars has been well studied [1], flight using passive sensors like cameras has relatively lagged behind. This is especially important given that small UAVs do not have the payload and power capabilities for carrying such sensors. Additonally, most of the modern research on UAVs has focussed on flying at altitudes with mostly open space [2]. Flying UAVs close to the ground through dense clutter [3] has been less explored.

Receding horizon control [4] is a classical deliberative scheme that has achieved much success in autonomous ground vehicles including five out of the six finalists of the DARPA Urban Challenge [5]. Figure 2 shows an illustration of receding horizon control on our UAV in motion capture. In receding horizon control, a pre-selected set of dynamically feasible trajectories of fixed length (the horizon), are evaluated on a cost map of the environment around the vehicle and the trajectory that avoids collision while making most progress towards a goal location is chosen. This trajectory is traversed for a bit and the process repeated again.

We demonstrate the *first* receding horizon control with monocular vision implementation on a UAV. Figure 1 shows our quadrotor evaluating a set of trajectories on the projected depth image obtained from monocular depth prediction and traversing the chosen one.

1, 2, 3, 7, 8 and 9 are with The Robotics Institute, Carnegie Mellon University, Pittsburgh, PA, USA {debadeep, kumarsha, samlzeng, daftry, hebert, dbagnell}@ri.cmu.edu

4 worked on the project while a student at Carnegie Mellon University, currently at NVIDIA Corporation, Santa Clara, CA, USA rupeshm@nvidia.com,

5 worked on the project while on an internship at Carnegie Mellon University, George Mason University, Fairfax, VA, USA, magcayaz@gmu.edu

6 worked on the project while on an internship at Carnegie Mellon University, Harvey Mudd College, Claremont, CA, USA, ceriksen@hmc.edu
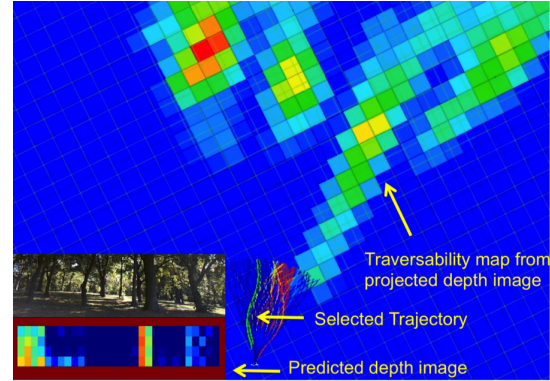
Fig. 1. Example of receding horizon with a quadrotor using monocular vision. The lower left images show the view from the front camera and the corresponding depth images from the monocular depth perception layer. The rest of the figure shows the overhead view of the quadrotor and the traversability map (built by projecting out the depth image) where red indicates higher obstacle density. The grid is $1x1$ $m^2$. The trajectories are evaluated on the projected depth image and the one with the least collision score (thick green) trajectory followed.
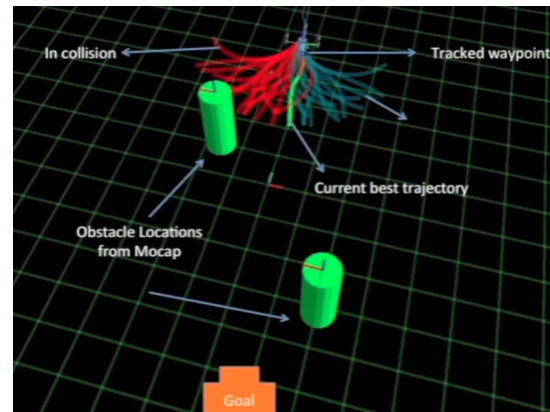


Fig. 2. Receding horizon control on UAV in motion capture. A library of 78 trajecories of length 5 m are evaluated to find the best collision-free trajectory. This is followed for some time and the process repeated.

This is motivated by our previous work [3] we used imitation learning to learn a purely reactive controller for flying a UAV using only monocular vision through dense clutter. While good obstacle avoidance behavior was obtained, there are certain limitations of a purely reactive layer that a more deliberative approach like receding horizon control can ameliorate. Reactive control is by definition myopic, i.e. it concerns itself with avoiding the obstacles closest to the vehicle. This can lead to it being easily stuck in cul-de-sacs. Since receding horizon control plans for longer horizons it achieves better plans and minimizes the chances of getting stuck [6]. Another limitation of pure reactive control is the difficulty to reach a goal location or direction. In a receding

horizon control scheme, trajectories are selected based on a score which is the sum of two terms: first, the collision score of traversing it and second, the heuristic cost of reaching the goal from the end of the trajectory. By weighting both these terms suitably, goal-directed behavior is realized while maintaining obstacle-avoidance capability. But it is to be noted that reactive control can be integrated with receding horizon for obtaining the best of both worlds in terms of collision avoidance behavior.

Receding horizon control needs three working components

- *A method to estimate depth*: This can be obtained from stereo vision [7], [8] or dense structure-from-motion (SfM) [9]. But these are not amenable for achieving higher speeds due to high computational expense. We note that in the presence of enough computation power, information from these techniques can be combined with monocular vision to improve overall perception. Biologists have found strong evidence that birds and insects use optical flow to navigate through dense clutter [10]. Optical flow has been used for autonomous flight of UAVs [11]. However, it tends to be very noisy and therefore difficult to directly derive a robust control principle from. Therefore we follow the same data driven principle as our previous work [3] and use local statistics of optical flow to serve as features in the monocular depth prediction module. This allows the learning algorithm to derive complex behaviors in a data driven fashion.
- *A method for relative pose estimation*: In order to track the chosen trajectory at every cycle, it is needed to estimate the relative state of the vehicle for the duration of the tracking. We developed and demonstrate a relative pose estimation system using a downward facing camera and a sonar to estimate height (Section II-E)
- *A method to deal with perception uncertainty*: Most planning schemes either assume that perception is perfect or make simplistic assumptions of uncertainty. We introduce the concept of making multiple, relevant yet diverse predictions for incorporating perception uncertainty into planning. The intuition is predicated on the observation that avoiding a small number of ghost obstacles is acceptable as far as true obstacles are not missed (high recall, low precision). The details are presented in Section II-D. We demonstrate in experiments the efficacy of this approach as compared to making only a single best prediction.

Our list of contributions in this work are:

- Budgeted near-optimal feature selection and fast non-linear regression for monocular depth prediction
- Real time relative vision-based pose estimation
- Multiple predictions to efficiently incorporate uncertainty in the planning stage.
- First complete receding horizon control implementation on a UAV.

Section II-A describes our hardware and software setup, while Section II-D and Section II-B describe the multiple predictions approach to handling uncertainty, and the budgeted depth prediction pipeline. The high frame-rate optical flow based pose estimation system developed for short-term

trajectory tracking is described in Section II-E. Section II-F describes the planning and control framework in detail, and Section III describes the outdoor experiments and results.

## II. APPROACH

Developing and testing all the integrated modules of receding horizon is very challenging. Therefore we developed a testing protocol where we assembled a rover in addition to a UAV to be able to test various modules separately. Here we describe the hardware platforms and software architecture of our system.

### A. Hardware and Software Overview

In this section we describe the hardware platforms used in our experiments.

In order to facilitate parallel testing of perception and planning modules while the UAV control systems were being developed, we assembled a ground rover (Figure 3) using off-the-shelf parts.
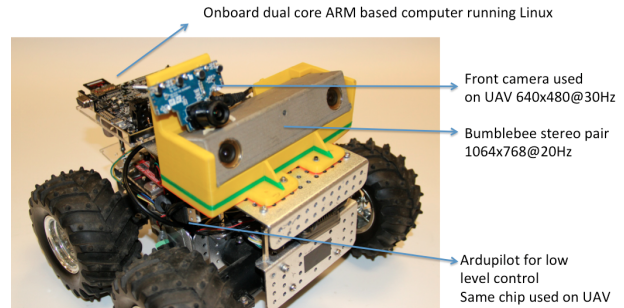


Onboard dual core ARM based computer running Linux

Front camera used on UAV 640x480@30Hz

Bumblebee stereo pair 1064x768@20Hz

Ardupilot for low level control Same chip used on UAV

Fig. 3. Rover assembled with the same control chips and perception software as UAV for rapid tandem development and validation of modules.

*1) Rover:* The rover (Figure 3) uses an Ardupilot microcontroller board [12] which takes in high level control commands from the planner and controls four motors to achieve the desired motion. The rover is skid-steered. Each motor also has integrated wheel encoders at 1633 counts per revolution. In order to get short-term relative pose for following a trajectory in receding-horizon control we developed an odometry package which runs on the Ardupilot.

Other than the low-level controllers all other aspects of the rover are kept exactly the same as the UAV to allow seamless transfer of software. For example the rover has a front facing PlayStation Eye camera which is also used as the front facing camera on the UAV.

A Bumblebee color stereo camera pair ($1024 \times 768$ at 20 fps) is rigidly mounted with respect to the front camera using a custom 3D printed fiber plastic encasing. This is used for collecting data with groundtruth depth values (Section II-B) and validation of planning (Section II-F). We calibrate the rigid body transform between the front camera and the left camera of the stereo pair. Stereo depth images and front camera images of the environment are recorded simultaneously while driving the rover around using a joystick. The depth images are then transformed to the front camera's coordinate system to provide groundtruth depth values for every pixel.

*2) UAV:* Figure 4 shows the quadrotor we use for our experiments. Figure 5 shows the schematic of the various modules that run onboard and offboard. The base chassis, motors and autopilot are assembled using the Arducopter kit [12]. Due to the excessive drift and noise of the IMU integrated in the Ardupilot unit, we added a Microstrain 3DM GX3 25 IMU which is used to aid real time pose estimation. There are two PlayStation Eye cameras: one facing downwards for real time pose estimation, one facing forward. The onboard processor is a quad-core ARM based computer which runs Ubuntu and ROS [13]. This unit runs the pose tracking and trajectory following modules. A sonar is used to estimate altitude. The image stream from the front facing camera is streamed to the base station where the depth prediction module processes it; the trajectory evaluation module then finds the best trajectory to follow to minimize probability of collision and transmits it to the onboard computer where the trajectory following module runs a pure pursuit controller to do trajectory tracking [14]. The resulting high level control commands are sent to the Ardupilot which sends low level control commands to the motor controllers to achieve the desired motion.
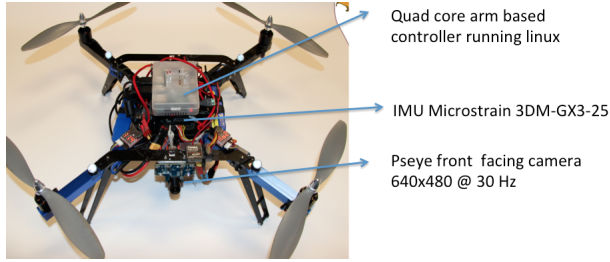


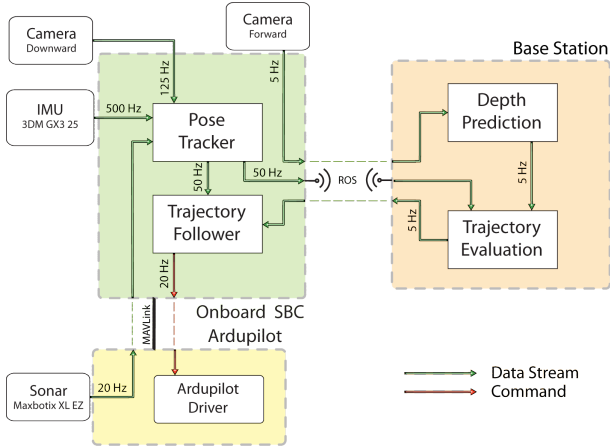Fig. 4. Quadrotor used as our development platform.



Fig. 5. Schematic diagram of hardware and software modules

## B. Monocular Depth Prediction

In this section we describe the depth prediction approach from monocular images, budgeted feature selection and the fast non-linear regression method used for regression.

In image is first gridded up into non-overlapping patches. We predict the depth in meters at every patch of the image (Figure 6 yellow box). For each patch we extract features which describe the patch, features which describe the full column containing the patch (Figure 6 green box) and features which describe the column of three times the patch width (Figure 6 red box), centered around the patch. The final feature vector for a patch is the concatenation of the feature vectors of all three regions. When a patch is seen by itself it is very hard to tell the relative depth with respect to the rest of the scene. But by adding the features of the surrounding area of the patch, more context is available to aid the predictor.



Fig. 6. The yellow box is an example patch, the green box is the column of the same width surrounding it, and the red box is the column of 3 times the patch width surrounding it. Features are extracted individually at the patch, and the columns and concatenated together to form the total feature representation of the patch.

*1) Description of features:* In this part we describe in brief the features used to represent the patch. We use the features as used in previous work on monocular imitation learning [3] for UAVs, which are partly inspired by the work of Saxena et al., [15]. In this case, instead of regressing directly from features of an image to a control action, we predict the depth at every patch which is then used by the planning module.

- *Optical flow*: We use the Farneback dense optical flow [16] implementation in OpenCV to compute for every patch the average, minimum and maximum optical flow values which are used as feature descriptors of that patch.
- *Radon Transform*: The radon transform captures strong edges in a patch [17].
- *Structure Tensor*: The structure tensor describes the local texture of a patch [18].
- *Laws' Masks*: These describe the texture intensities [19]. For details on radon transform, structure tensor and Laws' masks usage see [3].
- *Histogram of Oriented Gradients (HoG)*: This feature has been used widely in the computer vision community for capturing texture information for human pose estimation as well as object detection [20]. For each patch we compute the HoG feature in 9 orientation bins. We subdivide the patch into cells such that each resulting cell's side is at least half the width of the patch.
- *Tree feature*: We use the per pixel fast classifier by Li et al. [21] to train a supervised tree detector. Li et al. originally used this for real time hand detection in egocentric videos. For a given image patch we use this predictor to output the probability of each pixel being a tree. This information is then used as a feature for that patch.

*2) Data Collection:* RGB-D sensors like the Kinect, currently do not work outdoors. Since camera and calibrated nodding lidar setup is expensive and complicated we used a rigidly mounted Bumblebee stereo color camera and the PlayStation Eye camera for our outdoor data collection. This

setup was mounted on the rover as shown in Figure 3. We have collected data at 4 different locations with tree density varying from low to high, under varying illumination conditions and in both summer and winter conditions. Our corpus of imagery with stereo depth information is around 16000 images and growing. We plan on making this dataset publicly available in the near future.

*3) Fast Non-linear Prediction:* Due to harsh real-time constraints an accurate but fast predictor is needed. Recent linear regression implementations are very fast and can operate on millions of features in real time [22] but are limited in predictive performance by the inherent linearity assumption. In very recent work Agarwal et al. [23] develop fast iterative methods which use linear regression in the inner loop to obtain overall non-linear behavior. This leads to fast prediction times while obtaining much better accuracy. We implemented Algorithm 2 in [23] and found that it lowered the error by 10 % compared to just linear regression, while still allowing real time prediction.

### C. Budgeted Feature Selection

While there are many kinds of visual features that can be extracted, they need to be computed in real time. The faster the desired speed of the vehicle, the faster the perception and planning modules have to work to maintain safety. Compounded with the limited computational power available onboard a small UAV, this imposes a very challenging budget within which to make a prediction. Each kind of feature requires different time periods to extract, while contributing different amounts to the prediction accuracy. For example, radon transforms might take relatively less time to compute but contribute a lot to the prediction accuracy, while another feature might take more time but also contribute relatively less or vice versa. This problem is further complicated by the "grouping" effects where a particular feature's performance is affected by the presence or absence of other features.

Given a time budget, the naive but obvious solution is to enumerate all possible combinations of features which respect the budget and find the group of features which achieve the minimum loss. This is exponential in the number of available features. Instead we use the efficient approach developed by Hu et al. [24] to efficiently select the near-optimal set of features which meet the imposed budget constraints. Their approach uses a simple greedy algorithm that first whitens feature groups and then recursively chooses groups by the reduction in explained variance divided by the time to achieve that reduction. A more efficient variant of this with equivalent guarantees, chooses features by computing gradients to approximate the reduction in explained variance, eliminating the need to "try" all feature groups sequentially. For each specified time budget, the features selected by this procedure are within a constant factor of the optimal set of features which respect that time budget. Since this holds across all time budgets, this procedure provides a recursive way to generate feature sets across time steps.

Figure 7 shows the sequence of features that was selected by Hu et al.'s [24] feature selection procedure. For any given budget only the features on the left up to the specified time budget need to be computed.



Fig. 7. On the upper x-axis the sequence of features selected by Hu et al.'s method [24] and the lower x-axis shows the cumulative time taken for all features up to that point. The near-optimal sequence of features rapidly decrease the prediction error. For a given time budget, the sequence of features to the left of that time should be used.



Fig. 8. Depth prediction examples on real outdoor scenes. Closer obstacles are indicated by red.

### D. Multiple Predictions

The monocular depth estimates are often noisy and often inaccurate due the inherent challenging nature of the problem. A planning system must incorporate this uncertainty to achieve safe flight. Figure 9 illustrates the difficulty of trying to train a predictive method for building a perception system for any general collision avoidance problem. Figure 9 (left) shows a ground truth location of trees in the vicinity of an autonomous UAV flying through the forest. Figure 9 (middle) shows the location of the trees as predicted by the onboard perception system. In this prediction the trees on the left and far away in front are predicted correctly but the tree on the right is predicted close to the UAV. This will cause the UAV to dodge a ghost obstacle. While this is bad it is not fatal because the UAV will not crash but just have to make some spurious motions. But the prediction of trees in Figure 9 (right) is potentially fatal. Here the trees far away in front and on the right are correctly predicted where as the tree on the left originally close to the UAV is mis-predicted to be far away. This type of mistake will cause the UAV to

crash into an obstacle it does not know is there.



Fig. 9. Illustration of the complicated nature of the loss function for collision avoidance. (Left) Groundtruth tree locations (Middle) Bad prediction where a tree is predicted closer than it actually is located (Right) Fatal prediction where a tree close by is mis-predicted further away.

Ideally, a vision-based perception system should be trained to minimize loss functions which will penalize such fatal predictions more than other kind of predictions. But even writing down such a loss function is difficult. Therefore most monocular depth perception systems try to minimize easy to optimize surrogate loss functions like regularized $L_1$ or $L_2$ loss [15]. We try to reduce the probability of collision by generating multiple interpretations of the scene in front of the UAV to hedge against the risk of committing to a single interpretation of the scene which could be potentially fatal as illustrated above in Figure 9. We generate 3 interpretations of the scene and evaluate the trajectories in all 3 interpretations simultaneously. The trajectory which is least likely to collide on average in all interpretations is then chosen as the trajectory to traverse.

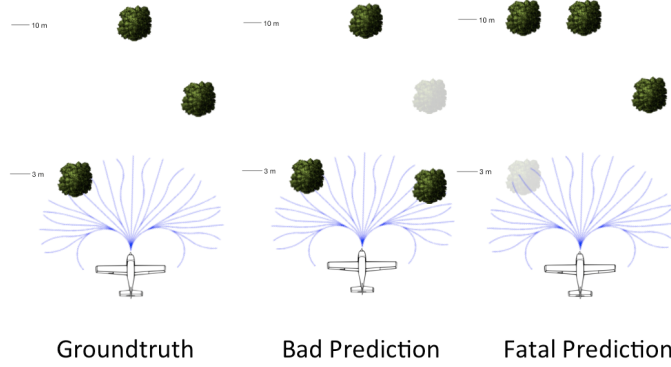One way of making multiple predictions is to just sample the posterior distribution of a learnt predictor. In order to truly capture the uncertainty of the predictor, a lot of interpretations have to be sampled and trajectories evaluated on each of them. A large number of samples will be from around the peaks of this distribution leading to wasted samples. This is not feasible given the real time constraints of the problem.

In previous work [25], we have developed techniques for predicting a budgeted number of interpretations of an environment with applications to manipulation, planning and control. Batra et al., [26] have also applied similar ideas to structured prediction problems in computer vision. These approaches try to come up with a small number of relevant but diverse interpretations of the scene so that at least one of them is correct. In this work, we adopt a similar philosophy and use the error profile of the fast non-linear regressor described in Section II-B to make two additional predictions: The non-linear regressor is first trained on a dataset of 14500 images and it's performance on a held-out dataset of 1500 images is evaluated. For each depth value predicted by it, the average over-prediction and under-prediction error is recorded. For example the predictor may say that an image patch is at 3 meters while on average whenever it says so, it is actually either, on average, at 4 meters or at 2.5 meters. We round each prediction depth to the nearest integer, and record the average over and under-predictions as in the above example in a look-up table (LUT). At test time the predictor

produces a depth map and the LUT is applied to this depth map, producing two additional depth maps: one for over-prediction error, and one for the under-prediction error.

Figure 10 shows an example in which making multiple predictions is clearly beneficial compared to the single best interpretation. We provide more experimental details and statistics in Section III.



Fig. 10. The scene at top is an example from the front camera of the UAV. On the left is shown the predicted traversability map (red is high cost, blue is low cost) resulting from a single interpretation of the scene. Here the UAV has selected the straight path (thick, green) which will make it collide with the tree right in front. While on the right the traversability map is constructed from multiple interpretations of the image, leading to the trajectory in the right being selected which will make the UAV avoid collision.

*E. Pose Estimation*

As discussed before, a relative pose-estimation system is needed to follow the trajectories chosen by the planning layer. We use a downward looking camera in conjunction with a sonar for determining relative pose. Looking forward to determine pose is ill-conditioned due to a lack of parallax as the camera faces the direction of motion. There are still significant challenges involved when looking down. Texture is often very self similar making it challenging for traditional feature based methods to be employed.

Dense methods such as DTAM [27] build a complete 3D model of the world that make tracking very easy, but require GPUs to run in real time. As a consequence, these can not be deployed onboard our platform. The other alternative is to use sparse mapping approaches like [28] that extract robustly identifiable interest points that are tracked over consecutive frames using a wide variety of feature descriptors, but fail due to issues mentioned above.

In receding horizon we do not need absolute pose with respect to some fixed world coordinate system. This is because we need to follow trajectories for short durations only. So as long as we have a relative, consistent pose estimation system for this duration (3 seconds in our implementation) we can successfully follow trajectories. Thus by using only differential motion between consecutive frames and not consciously refining pose information, we can get away with small drift over time.

This motivated the use of a variant of a simple algorithm that has been presented quite often, most recently in [29]. This approach involves using a Kanade-Lucas-Tomasi (KLT) tracker [30] to detect where each pixel in a grid of pixels

Fig. 11. The overall flow of data and control commands between various modules. The pure pursuit trajectory follower and low level control loops (red boxes) are shown in greater detail at the bottom.

moves over consecutive frames, and estimating the mean flow from these after rejecting outliers. The KLT tracker operates on the principle of assuming a small consistent flow from a patch around a feature point in the first image to another in the second image. It does this by performing iterative Newton-Raphson refinement steps minimizing the error in image intensities for each pixel over the patch weighted by the local second derivative around the interest point. This estimate of flow however tries to find the best planar displacement between the two patches, and does not take into account out-of-plane rotations, for instance, due to motion of the camera. Camera ego-motion is compensated using motion information from the IMU. Finally the scene depth is estimated from a sonar. We obtain instantaneous relative velocity between the camera and ground which we integrate over time to get position.

This process is computationally inexpensive, and can thus be run at very high frame rates. Higher frame rates lead to smaller displacements between pairs of images, which in turn makes tracking easier to compute and more accurate since over short displacements first order approximations to nonlinearities have less error, thus driving down the overall error in the system.

*1) Implementation Details:* We obtain the angular velocity of the camera using the output of the IMU driver that provides us orientation and velocity updates at upto 500Hz. The optical flow of the central patch is determined by using OpenCV's pyramidal Lucas-Kanade tracker implementation.

For robustness to outliers, we evaluate flow for a small discrete grid around the central patch and take the mean of the unrotated optical flow vectors as the final flow vector. We also do an outlier detection step by comparing the standard deviation of the flow vectors obtained for every pixel on the grid to a specific threshold. Whenever the variance of the flow is high, we do not calculate the mean flow velocity, and instead decay the previous velocity estimate by a constant factor. However, such situations rarely occur due to the very high frame rate of the camera feed we use that ensures that only a small change in the flow field occurs between consecutive images.

*2) Performance vs Ground Truth:* To evaluate the performance of the differential flow based tracker, we took the

camera-IMU quadrotor setup to the motion capture lab and compared the tracking performance. We held the quadrotor in our hand and walked over artificial camouflage texture. The resulting tracks are as shown in figure 12



Fig. 12. Comparison of the differential flow tracker performance vs ground truth in MOCAP. Red tracks are the trajectories in MOCAP, blue are those determined by the algorithm. Note that the formulation of the receding horizon setup is such that mistakes made in following a specific trajectory are forgiven up to an extent since we replan every few seconds.

### F. Planning and Control

Figure 11 shows the overall flow of data and control commands in the architecture. The front facing camera video stream is fed to the perception module which predicts the depth of every pixel in a frame, projects it to a point cloud representation and sends it to the receding horizon control module. A trajectory library of 78 trajectories of length 5 meters is budgeted and picked from a much larger library of 2401 trajectories using the maximum dispersion algorithm by Green et al. [31]. This is a greedy procedure for selecting trajectories, one at a time, so that each subsequent trajectory spans maximum area between it and the rest of the trajectories. The receding horizon module maintains a score value for every point in the point cloud. The score of a point decays exponentially the longer it exists. After some time when it drops below a user set threshold the point is deleted. The rate of decay is specified by setting the time constant of the decaying function. This fading memory representation of the local scene layout has two advantages: 1) It prevents collisions caused by narrow field-of-view issues where the quadrotor forgets that it has just avoided a tree, sees the next tree and dodges sideways, crashing into the just avoided tree. 2) It allows emergency backtracking maneuvers to be safely executed if required, since there is some local memory of the obstacles it has just passed.

Our system can accept any goal direction as input and ensures that the vehicle makes progress towards the goal while avoiding obstacles along the way. The score for each trajectory is the sum of three terms: 1) A sphere of the same radius as the quadrotor is convolved along a trajectory



Fig. 13. Instances of failure of the pose tracking system over challenging surfaces. Note the absence of texture in these 320x240 images. The figure shows the flow tracks corresponding to the points on the grid. Red tracks show the uncorrected optical flow, while the green tracks (superimposed) show the flow vectors 'unrotated' using the IMU.

and the score of each point in collision is added up. The higher this term relative to other trajectories, the higher the likelihood of this trajectory being in collision. 2) A term which penalizes a trajectory whose end *direction* deviates from goal direction. This is weighted by a user specified parameter. This term induces goal directed behavior and is tuned to ensure that the planner always avoids obstacles as a first priority. 3) A term which penalizes a trajectory for deviating in *translation* from the goal direction.

The pure pursuit controller module (Figure 11) takes in the coordinates of the trajectory to follow and the current pose of the vehicle from the optical flow-based pose estimation system (Section II-E). We use a pure pursuit strategy [14] to successfully track it. Specifically, this involves finding the closest point on the trajectory from the robot's current estimated position and setting the target waypoint to be a certain fixed lookahead distance further along the trajectory. The lookahead distance can be tuned to obtain the desired smoothness while following the trajectory; concretely, a larger lookahead distance leads to smoother motions, at the cost of not following the trajectory exactly. Using the pose updates provided by the pose estimation module, we head towards this moving waypoint using a generic PD controller. Since the receding horizon control module continuously replans (at 5 hz) based on the image data provided by the front facing camera, we can choose to follow arbitrary lengths along a particular trajectory before switching over to the latest chosen one.

*1) Validation of Modules:* We validated each module separately as well as in tandem with other modules where each validation was progressively integrated with other modules. This helped reveal bugs and instabilities in the system.

- *Trajectory Evaluation and Pure Pursuit Validation with Stereo Data on Rover*: We tested the trajectory evaluation and pure pursuit control module by running the entire pipeline (other than monocular depth prediction) with stereo depth images on the rover. Figure 14.
- *Trajectory Evaluation and Pure Pursuit Validation with Monocular Depth on Rover*: This test is the same as above but instead of using depth images from stereo we used the perception system's depth prediction. This allowed us to tune the parameters for scoring trajectories in the receding horizon module to head towards goal without colliding with obstacles.
- *Trajectory Evaluation and Pure Pursuit Validation with Known Obstacles in Motion Capture on UAV*: While testing of modules progressed on the rover we assembled and developed the pose estimation module (Section II-E) for the UAV. We tested this module in a motion capture lab where the position of the UAV as well of the obstacles was known and updated at 120 Hz. (See Figure 2)
- *Trajectory Evaluation and Pure Pursuit Validation with Hardware-in-the-Loop (HWIL)*: In this test we ran the UAV in an open field, fooled the receding horizon module to think it was in the midst of a point cloud and ran the whole system (except perception) to validate planning and control modules. Figure 15 shows an example from this setup.
- *Whole System*: After validating each module following

the evaluation protocol described above, we ran the whole system end-to-end. Figure 1 shows an example scene of the quadrotor in full autonomous mode avoiding trees outdoors. We detail the results of collision avoidance in Section III.
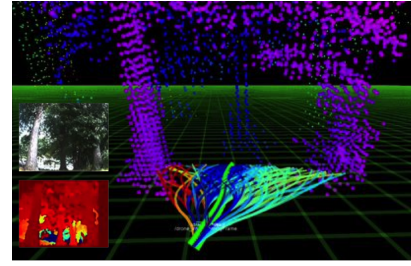


Fig. 14. Receding horizon control validation with rover using depth images from stereo. The bright green trajectory is the currently selected trajectory to follow. Red trajectories indicate that they are more likely to be in collision.
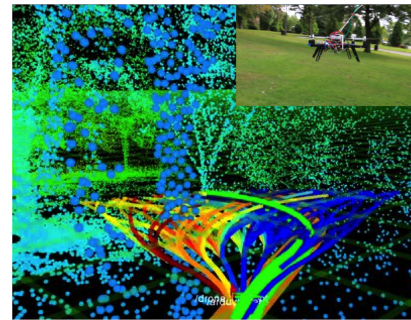


Fig. 15. Hardware-in-the-loop testing with UAV in open field. The receding horizon module was fooled into thinking that it was in the midst of a real world point cloud while it planned and executed its way through it. This allowed us to validate planning and control without endangering the UAV.

## III. EXPERIMENTS

In this section we analyze the performance of our proposed deliberative approach for autonomous navigation of a UAV in cluttered natural environments using only monocular vision. All the experiments were conducted in a densely cluttered forest area, while restraining the drone through a light-weight tether.

Quantitatively, we evaluate the performance of our system by observing the average distance flown autonomously by the UAV over several runs, before an intervention. An intervention, in this context, is defined as the point at which the pilot needs to overwrite the commands generated by our control system so as to prevent the drone from crashing. Experiments were performed using both our proposed multiple predictions approach and single best prediction, and the corresponding comparison has been shown in Fig. 16. Both tests were performed in regions of high and low clutter density (approx. 1 tree per 6x6 $m^2$ and 12x12 $m^2$, respectively). It can be observed that multiple predictions results in significantly better performance. In particular, the drone was able to fly autonomously without crashing over a 137 m distance for low density regions. The difference is even higher in case of high-density regions where committing to a single prediction can be even more fatal.

Further, we evaluate the success rate for avoiding the large and small trees using our proposed approach. The results

| | Multiple Predictions | Single Prediction |
|---|---|---|
| Total Distance | 1020 m | 1010 m |
| Large Trees Avoided | 93.1 % | 84.8 % |
| Small Trees Avoided | 98.6 % | 95.9 % |
| Overall Accuracy | **96.6 %** | **92.5 %** |

TABLE I

SUCCESS RATE OF AVOIDING TREES.

have been tabulated in Table I. We are successfully able to avoid 96 % of all trees over a total covered distance of more than 1 km. We extend our evaluation to qualitatively assess, and understand the failure cases responsible for the above results (See Fig. 16). The type of failures are broken down by the type of obstacle the drone failed to avoid, or whether the obstacle was not in the field-of-view (FOV). Overall, 39 % of the failures were due to large trees and 33 % on hard to perceive obstacles like branches and leaves. As expected, the narrow FOV is now the least contributor to the failure cases as compared to a more reactive control strategy [3]. This is intuitive, since the reactive control is myopic in nature and our deliberate approach helps overcome this problem as described in the previous sections. Figure 17 shows some typical intervention examples.
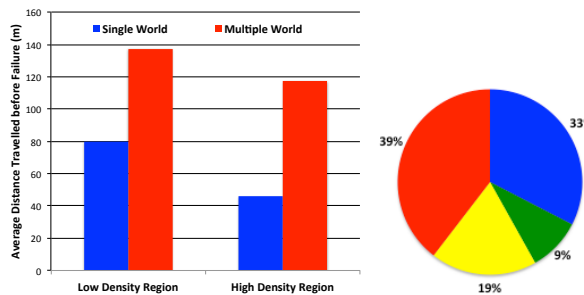


Fig. 16. (a) Average distance flown by the drone before a failure. (b) Percentage of failure for each type. Red: Large Trees, Yellow: Thin Trees, Blue: Foliage, Green: Narrow FOV.



Fig. 17. Examples of interventions: (Left) Bright trees saturated by sunlight from behind (Second from left) Thick foliage (Third from left) Thin trees (Right) Flare from direct sunlight. Camera/lens with higher dynamic range and more data of rare classes should improve performance.

## IV. CONCLUSION

While we have obtained promising results from the current approach, a number of challenges remain. For example, better handling of sudden strong wind disturbances and control schemes for better leveraging the full dynamic envelope of the vehicle. In ongoing work we are moving towards complete onboard computing of all modules to reduce latency. We can leverage other sensing modes like sparse, but more accurate depth estimation from stereo, which can be used as "anchor" points to improve dense monocular depth estimation. Similarly low power, light weight lidars can be actively foveated to high probability obstacle regions to reduce false positives and get exact depth.

Another central future effort is to integrate the purely reactive [3] approach with the deliberative scheme detailed here, for better performance.

## REFERENCES

[1] S. Scherer, S. Singh, L. J. Chamberlain, and M. Elgersma, "Flying fast and low among obstacles: Methodology and experiments," *IJRR*, 2008.

[2] D. Dey, C. Geyer, S. Singh, and M. Digioia, "A cascaded method to detect aircraft in video imagery," *IJRR*, 2011.

[3] S. Ross, N. Melik-Barkhudarov, K. S. Shankar, A. Wendel, D. Dey, J. A. Bagnell, and M. Hebert, "Learning monocular reactive uav control in cluttered natural environments," in *ICRA*, 2013.

[4] A. Kelly *et al.*, "Toward reliable off road autonomous vehicles operating in challenging environments," *IJRR*, 2006.

[5] M. Buehler, K. Iagnemma, and S. Singh, "Special issue on the 2007 darpa urban challenge, part i, ii, iii," *JFR*, 2008.

[6] R. Knepper and M. Mason, "Path diversity is only part of the problem," in *ICRA*, May 2009.

[7] K. Schmid, P. Lutz, T. Tomić, E. Mair, and H. Hirschmüller, "Autonomous vision-based micro air vehicle for indoor and outdoor navigation," *JFR*, 2014.

[8] L. Matthies, R. Brockers, Y. Kuwata, and S. Weiss, "Stereo vision-based obstacle avoidance for micro air vehicles using disparity space," in *ICRA*, 2014.

[9] A. Wendel, M. Maurer, G. Graber, T. Pock, and H. Bischof, "Dense reconstruction on-the-fly," in *CVPR*, 2012.

[10] M. V. Srinivasan, "Visual control of navigation in insects and its relevance for robotics," *Current opinion in neurobiology*, 2011.

[11] A. Beyeler, J.-C. Zufferey, and D. Floreano, "Vision-based control of near-obstacle flight," *Autonomous robots*, 2009.

[12] [Online]. Available: http://dev.ardupilot.com

[13] M. Quigley, K. Conley, B. P. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, and A. Y. Ng, "Ros: an open-source robot operating system," in *ICRA Workshop on Open Source Software*, 2009.

[14] R. C. Coulter, "Implementation of the pure pursuit path tracking algorithm," DTIC Document, Tech. Rep., 1992.

[15] A. Saxena, S. H. Chung, and A. Y. Ng, "Learning depth from single monocular images," in *NIPS*, 2005.

[16] G. Farnebäck, "Two-frame motion estimation based on polynomial expansion," in *Image Analysis*, 2003.

[17] S. Helgason, "Support of radon transforms," *Advances in Mathematics*, 1980.

[18] C. Harris and M. Stephens, "A combined corner and edge detector." in *Alvey vision conference*, 1988.

[19] E. R. Davies, *Machine vision: theory, algorithms, practicalities*, 2004.

[20] N. Dalal, B. Triggs, and C. Schmid, "Human detection using oriented histograms of flow and appearance," in *ECCV*, 2006.

[21] C. Li and K. M. Kitani, "Pixel-level hand detection in ego-centric videos," in *CVPR*, 2013.

[22] J. Langford, L. Li, and A. Strehl, "Vowpal Wabbit," 2007.

[23] A. Agarwal, S. M. Kakade, N. Karampatziakis, L. Song, and G. Valiant, "Least squares revisited: Scalable approaches for multi-class prediction," *arXiv preprint arXiv:1310.1949*, 2013.

[24] H. Hu, A. Grubb, J. A. Bagnell, and M. Hebert, "Efficient feature group sequencing for anytime linear prediction," *arXiv:1409.5495*, 2014.

[25] D. Dey, T. Y. Liu, M. Hebert, and J. A. D. Bagnell, "Contextual sequence optimization with application to control library optimization," in *RSS*, 2012.

[26] D. Batra, P. Yadollahpour, A. Guzman-Rivera, and G. Shakhnarovich, "Diverse m-best solutions in markov random fields," in *Computer Vision–ECCV 2012*. Springer, 2012, pp. 1–16.

[27] R. A. Newcombe, S. J. Lovegrove, and A. J. Davison, "Dtam: Dense tracking and mapping in real-time," in *ICCV*, 2011.

[28] G. Klein and D. Murray, "Parallel tracking and mapping for small ar workspaces," in *ISMAR*, 2007.

[29] D. Honegger, L. Meier, P. Tanskanen, and M. Pollefeys, "An open source and open hardware embedded metric optical flow cmos camera for indoor and outdoor applications," in *ICRA*, 2013.

[30] C. Tomasi and T. Kanade, *Detection and tracking of point features*. School of Computer Science, Carnegie Mellon Univ., 1991.

[31] C. Green and A. Kelly, "Optimal sampling in the space of paths: Preliminary results," Robotics Institute, Pittsburgh, PA, Tech. Rep. CMU-RI-TR-06-51, November 2006.