

Article

# Development of an Autonomous Unmanned Aerial Manipulator Based on a Real-Time Oriented-Object Detection Method

Shijie Lin <sup>1,†,ID</sup>, Jinwang Wang <sup>1,†</sup>, Rui Peng <sup>1</sup> and Wen Yang <sup>1,2,\*ID</sup>

<sup>1</sup> School of Electronic Information, Wuhan University, Wuhan 430072, China; linshijie@whu.edu.cn(S.L.); jwwangchn@whu.edu.cn(J.W.); pengrui@whu.edu.cn(R.P.)

<sup>2</sup> State Key Laboratory for Information Engineering in Surveying, Mapping and Remote Sensing (LIESMARS), Wuhan 430079, China

\* Correspondence: yangwen@whu.edu.cn(W.Y.)

† These authors contributed equally to this work.

Version March 6, 2019 submitted to Sensors

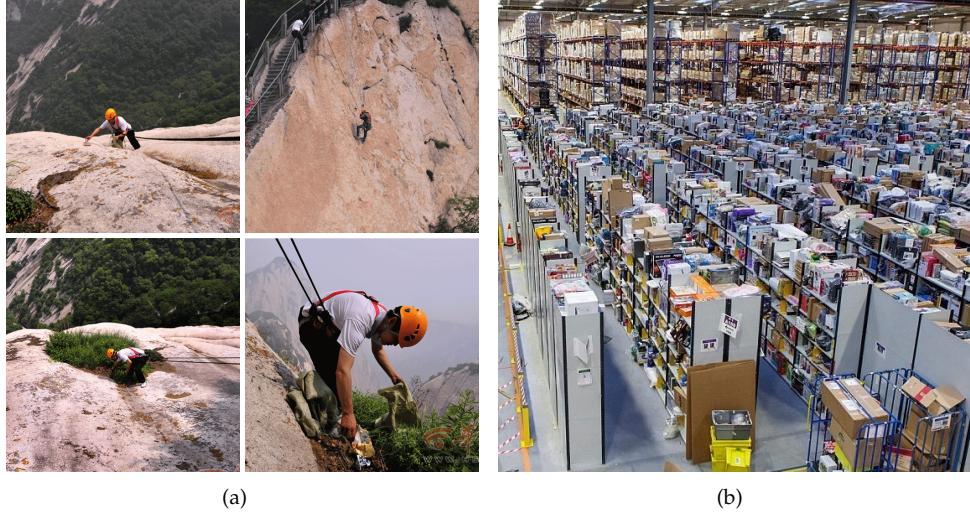
**Abstract:** Autonomous Unmanned Aerial Manipulators (UAMs) have shown promising potentials in mobile 3-dimension grasping applications, but they still suffer from some difficulties impeding their board applications, such as target detection and indoor positioning. For the autonomous grasping mission, the UAMs need ability to recognize the objects and grasp them. Considering the efficiency and precision, we present a novel oriented-object detection method called Rotation-SqueezeDet. This method can run on embedded-platforms in near real-time. Besides, this method can give the oriented bounding box of object in images to enable a rotation-aware grasping. Based on this method, a low-cost UAM platform form is designed and built. We give the formulation, positioning, control, and planning of the whole UAM system. All the mechanical designs are fully provided as open-source hardware for the reuse by the community. Finally, the effectiveness of the proposed scheme is validated in multiple experimental trials, highlighting its applicability of autonomous aerial rotational grasping in GPS-denied environments. We believe this system can be deployed to many potential workplaces which need UAM to accomplish difficult manipulation tasks.

**Keywords:** aerial manipulation; aerial system; deep learning

## 1. Introduction

During the past several decades, Unmanned aerial vehicles (UAVs) have shown surprising potentials among numerous applications[1,2]. Mounting with different kinds of sensors, the UAVs can extend the sensing range and transforming a static sensing task into a mobile sensing task. The UAVs' advantages of freely using the 3-dimension (3D) space bring possibilities to many traditional tasks like manipulation tasks in a warehouse or a factory.

Unmanned aerial manipulators (UAMs) are known as one specific type of UAVs equipped with one or multiple robotic arms and have attracted a lot of research interests in recent years [3]. One main advantage of UAM is that it shows promising potentials to transform passive sensing missions into mobile 3D interactive missions like grasping [4] and assembling [5]. The capabilities like aerial maneuvering and hovering make it possible for UAM to accomplish many kinds of mission that are difficult for human workers, such as the costly grasping bottles on the cliff shown in Fig. 1(a). Furthermore, there are many places like the Amazon warehouse shown in Fig. 1(b) can be the potential places to deploy UAM system for autonomous picking and placing. To fully utilize the space in a warehouse, goods can be placed at a relatively high place difficult for human workers to get. At this time, the UAM can give a lot of help. In the future warehouse, there will be many UAMs flying



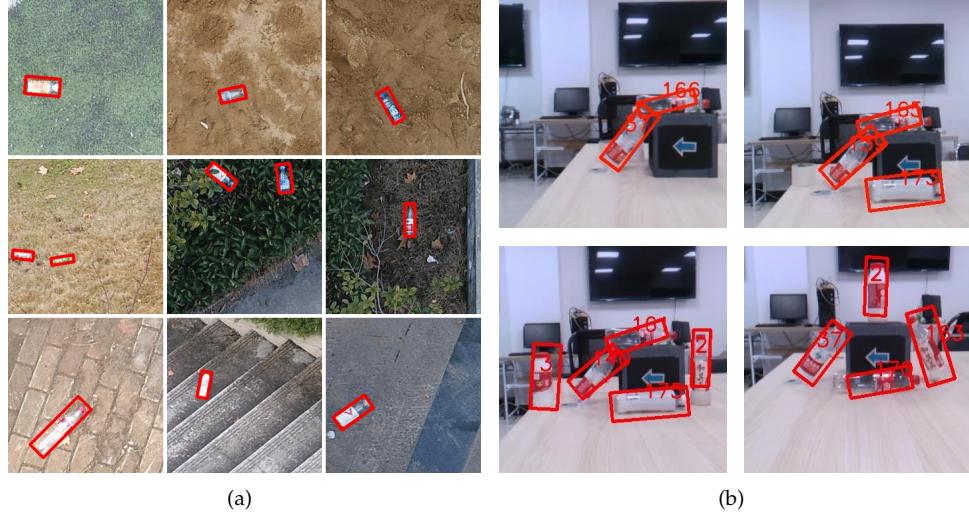
**Figure 1.** (a) Workers are carrying out dangerous garbage disposal work on the cliff. (b) The Amazon warehouse is filled with many sorts of goods.

31 simultaneously to sorting the goods. Besides, places like the Fukushima nuclear power plant in Japan  
 32 that filled will obstacles and hard for human workers to work inside can be the ideal places to deploy  
 33 UAM. Potential applications of UAMs can easily come up from people's imagination and surely not  
 34 being limited to the situations mentioned above that is the reason why we believe UAMs have a bright  
 35 future.

36 For the autonomous aerial manipulation missions, building a controllable system is always  
 37 the first step. However, many factors can influence the stability of the overall system, such as the  
 38 change of Center of Gravity (CoG) generated by the movements of the robotic arm, the reaction force  
 39 produced by the robotic arm, the complex aerodynamics effects. Many efforts have been done to  
 40 reduce these effects[6–10]. A comprehensive dynamic model of hexacopter with a robotic arm has  
 41 been built in [6], analyzing the effects of CoG change and mass distributions. The kinematics and  
 42 dynamics models of quadrotor with n Degree of Freedom (DoF) robotic arm have been formulated in  
 43 [7]. A Variable Parameter Integral Backstepping (VPIB) controller proposed in [8] can control a UAM  
 44 with better performance than PID controller. With a movable compensation mechanism, a multilayer  
 45 architecture controller compensates the internal and external effects layer by layer to control the UAM  
 46 was presented [9]. To suppress the torque generated by the movements of the robotic arm, a novel  
 47 mechanism with a simplified model has been adopted in [10].

48 Usually, human operators are unable to accurately control the UAM due to the point of view  
 49 change and data transmission delay. Also, for the monocular camera, the lack of depth perception  
 50 in teleoperation making the operators hard to decide whether the object is appropriately put into  
 51 the claw of robotic arm. Such inaccuracy causes the robotic arm to be difficult to align and grasp  
 52 the objects. Therefore it is better the UAM can work without relying on external control by human  
 53 operators. However, if there is no external control by human operators, the UAM needs to recognize  
 54 the objects by themselves. Meaning we need to give UAM the perception ability. While currently only  
 55 a few works[11–13] considered robust visual perception aided autonomous grasping. An Image-Based  
 56 Visual Servo (IBVS) was implemented in [11] to help to locate the position of the targeted object.  
 57 Feature models are used in [12] to find the position of known targets. In [11], correlation filters were  
 58 adopted to track targets.

59 In recent years, end-to-end object detection algorithms like [14], [15] and [16] have shown  
 60 surprising results. However, attempts directly apply these methods in the UAM system usually  
 61 end with failure or unsatisfactory robustness. The reason lies in the objects, like bottles, in UAV  
 62 perspective often appeared with arbitrary poses as shown in Fig. 2(a), since the camera is tight-coupled



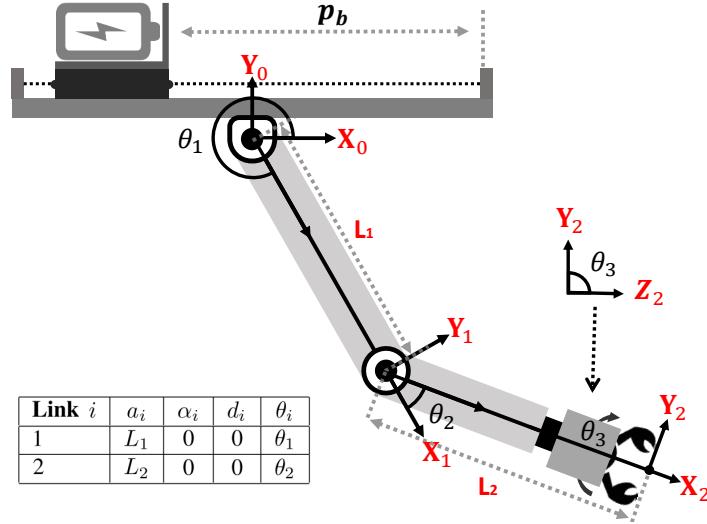
**Figure 2.** (a) Sample images in UAV-DB dataset with plastic bottle groundtruth labeled by red box. (b) Detection results of Rotation-SqueezeDet with 2D box and rotation angle of multiple objects at the same time. Oriented angles of bottles are shown in the images.

and rotated with the UAV body. Moreover, most grasping missions need a target's rotation angle for not touching the target when the end-effector is approaching. Furthermore, the performance like robustness and efficiency of the detection algorithms play an important role in grasping missions since it can bring more mobility to the UAM. However, current common oriented bounding boxes based methods presented in [17] and [18] are designed for large-scale objects detection like detected a car in a large satellite image. Since the cars only take a few pixels in a large satellite image, they need a much deeper convolution network to improve their accuracy. And this deep network cannot run in real-time performance even with an NVIDIA GTX1080 graphic card. Considering the size and weight requirement of a drone platform, the NVIDIA Jetson TX2 is the current optimal choice for mobile deep learning platforms. But methods like [17,18] are still unable to run onboard since both require more than 11G GPU memory to load a deep convolution network and Jetson TX2 only has 8G GPU memory, let alone the real-time performance. Without a real-time performance, the grasping efficiency has been greatly reduced because the UAM need to fly much slower for not passing the targets.

In order to solve the problems mentioned above, we propose Rotation-SqueezeDet which can regress rotation angle and position in 2-dimension (2D) image in near real-time. Unlike the common horizontal bounding box descriptor  $(c_x, c_y, w, h)$ , where  $(c_x, c_y)$  is the center location,  $w$  and  $h$  are the width and height of the bounding box, respectively, Rotation-SqueezeDet introduces a new  $\theta$  term, and thus uses  $(c_x, c_y, w, h, \theta)$  to describe object position and rotation angle in the 2D image. This not only makes the detection more robust since the bounding box included fewer background, but also provides the rotation angle of the target. By using an Intel RealSense D435 depth camera, the relative 3D distance of targets can be measured in the point cloud generated by registered depth image once the target is detected. Hence, a rotation-aware grasping for autonomous grasping is possible. A glimpse of detection results is shown in Fig. 2(b).

Now we want to deploy this detection in the UAM for autonomous grasping. However, the system mentioned in previous paragraphs have only considered neither the control problem or system compensate problem but an overall applicable indoor system. Usually, they test system in an outdoor environment with strong GPS singles for localization or use the expensive external visual caption system like VICON<sup>1</sup> to provide state estimation of drone. Both of these two ways are hard to achieve in

<sup>1</sup> [www.vicon.com](http://www.vicon.com)



**Figure 3.** Coornidates system and D-H parameter of the robotic arm.

91 a real industrial environment. So we adopted the lightweight SLAM system VINS[19] for indoor state  
 92 estimation, give our system the ability to fly indoor with only an inexpensive camera. Furthermore,  
 93 based on the Rotation-SqueezeDet, we design and implement a complete indoor workable UAM  
 94 system.

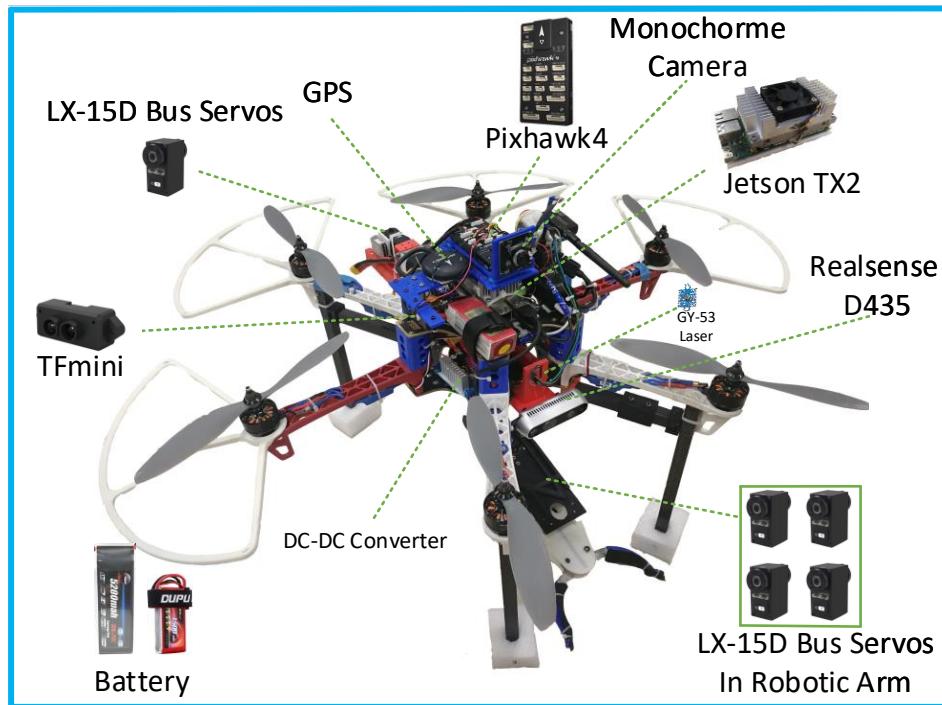
95 The main contributions of this paper can be summarized in two folds. First, the  
 96 Rotation-SqueezeDet method is proposed. This method can run on Jetson TX2 in near real-time  
 97 and enable successful rotation-aware grasping. What's more, we believe that this method will be  
 98 suitable in not only the aerial grasping but also a more general missions. Second, we have designed  
 99 and assembled the whole UAM platform, fully tested this system in an indoor environment. The  
 100 system is affordable since it costs less than \$2300 USD and can fly without relying on expensive  
 101 visual motion caption system. We believe this system can be adopted to many potential workplaces  
 102 like the Amazon warehouse which needs UAM to accomplish difficult manipulation tasks. All  
 103 mechanical structures are provided as open-hardware for reuse by the community. Link: <https://github.com/eleboss/UAMmech>

105 The rest of this paper is organized as follows. Section 2 describes the design and details of the  
 106 overall system. Section 3 describes the formulation, control, and planning of the UAM system. Section  
 107 4 describes the complete vision system including Rotation-SqueezeDet. Section 5 experimentally  
 108 demonstrates the system including autonomous grasping and vision system performance in an open  
 109 dataset. Finally, the conclusion and future work are presented in Section 6.

## 110 2. System Description

### 111 2.1. Notation

112 The East, North, and Up (ENU) coordinate system are used as a world-fixed inertial frame  
 113 corresponding to  $\{x_w, y_w, z_w\}$ . Following the definition of well known Denavit-Hartenberg (D-H)  
 114 parameters [20], the link frame of Link  $i$  is defined as the  $\{x_i, y_i, z_i\}$ ,  $i = 0$  is the fixed arm frame.  
 115 The definition of the link frame is detailed in Fig. 3 and the table in the bottom left gives the D-H  
 116 parameters of the robotic arm. The body frame is assumed to be the geometrical center of the UAM  
 117 denoted as  $\{x_b, y_b, z_b\}$ .  $G_x, G_y, G_z$  indicate the CoG of the UAM in body frame  $\{x_b, y_b, z_b\}$ .  $(\phi, \theta, \psi)$   
 118 indicate the roll-pitch-yaw Euler angles.  $\{x_t, y_t, z_t\}$  define the detected targets in the RealSense D435  
 119 camera frame.



**Figure 4.** Hardware components of the UAM

## 120 2.2. Hardware

121 In this work, a modified DJI hexacopter frame F550<sup>2</sup> is adopted. The hardware components  
 122 are shown in Fig. 4. The propulsion module of the UAM is composed of Sunnysky<sup>3</sup> x3108s motor,  
 123 Hobbywing<sup>4</sup> platinum Electronic Speed Controller (ESC), and 10 inches propeller. The Pixhawk4<sup>5</sup>  
 124 autopilot with PX4<sup>6</sup> V1.8.0 flight stack and NVIDIA Jetson TX2<sup>7</sup> are placed at the top of the UAM as  
 125 the main computing devices. A global shutter monochrome camera<sup>8</sup> is tight-coupled with Pixhawk4<sup>5</sup>  
 126 by using a 3D-printed anti-vibration damping plate. A Benewake TFmini<sup>9</sup> laser rangefinder is chosen  
 127 for being cheap, lightweight and with up to 12m maximum detection range, mounted downward  
 128 facing to provide altitude feedback. An Intel RealSense D435<sup>10</sup> camera is mounted at the middle of the  
 129 drone facing forward to find the targets.

130 We build a displacement compensation system (DCS) which can move counterweight to align the  
 131 CoG and thus improves the stability of the total system. The DCS is mounted in the middle of UAM  
 132 and made by 3D printed PLA material, including tow rails, a slide table and a bus servo to provide  
 133 drive force. The LEBOT<sup>11</sup> LX-15D serial bus servos are chosen for being budget-friendly, lightweight,  
 134 and having multiple extra structures to facilitate the installation. Another significant advantage of  
 135 the bus servo is it can largely reduce the wiring complexity and control difficulty. So we can only  
 136 use one serial port in Jetson TX2 to control all servos. While the LX-15D servo can only provide 240°

<sup>2</sup> [www.dji.com/flame-wheel-arf/feature](http://www.dji.com/flame-wheel-arf/feature)

<sup>3</sup> [www.rcsunnysky.com](http://www.rcsunnysky.com)

<sup>4</sup> [www.hobbywing.com](http://www.hobbywing.com)

<sup>5</sup> [www.holybro.com/product/55](http://www.holybro.com/product/55)

<sup>6</sup> <https://github.com/PX4/Firmware>

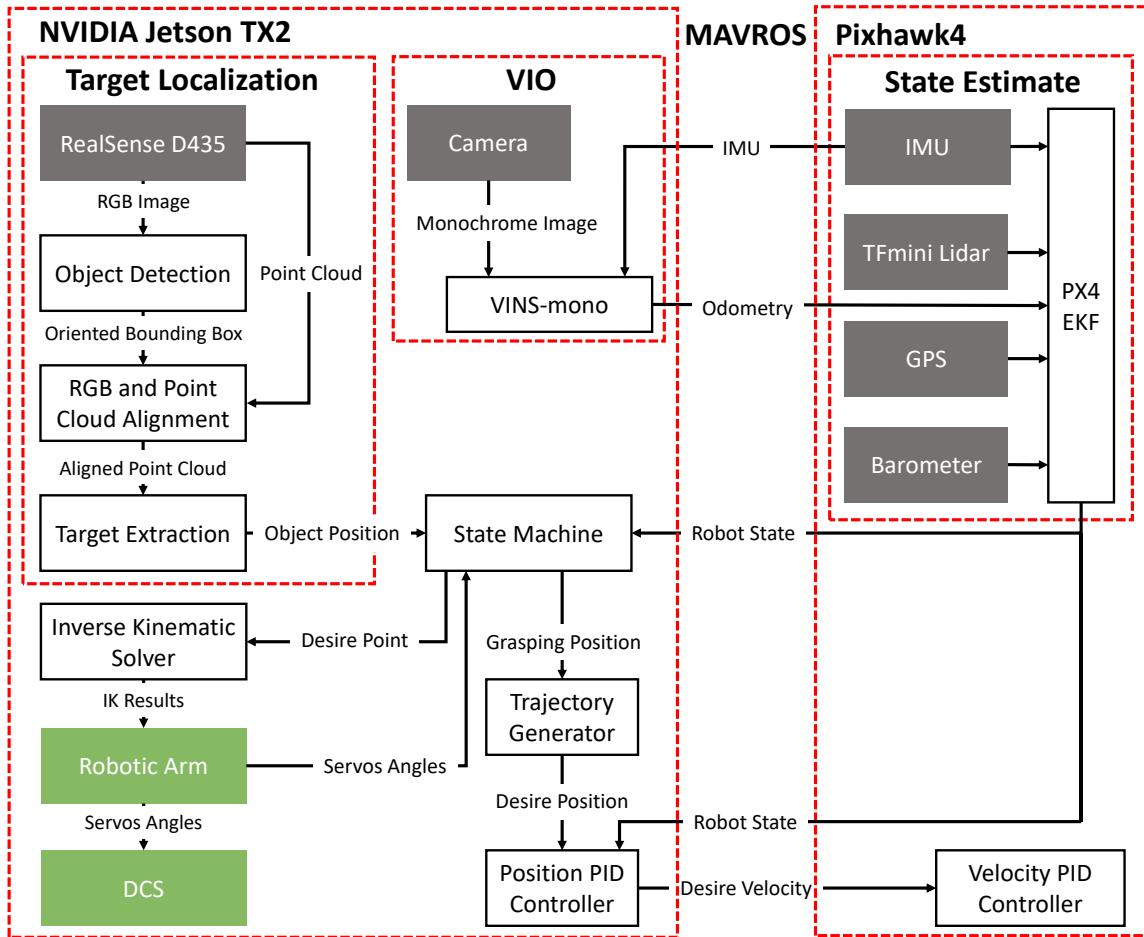
<sup>7</sup> [developer.nvidia.com/embedded/buy/jetson-tx2](http://developer.nvidia.com/embedded/buy/jetson-tx2)

<sup>8</sup> [www.jinyandianzi.com](http://www.jinyandianzi.com)

<sup>9</sup> [www.benewake.com/tfmini.html](http://www.benewake.com/tfmini.html)

<sup>10</sup> [www.realsense.intel.com/stereo](http://www.realsense.intel.com/stereo)

<sup>11</sup> [www.lobot-robot.com](http://www.lobot-robot.com)



**Figure 5.** An overview of the software architecture.

137 feedback, we use a short-range time of flight (ToF) laser rangefinder GY-53 to provide the battery  
 138 position feedback.

139 A 5200mAh 4S-35C battery weighted 0.525kg is used for providing enough power to the  
 140 propulsion system and as a counterweight for the DCS. And this battery can sustain the flight time  
 141 around 10min. Another 1500mAh 3S-30C weight 0.138kg battery is used for providing power to the  
 142 robotic arm and computing facilitates.

143 Drive forces of robotic arm are also provided by LX-15D servos. The robotic arm has 3-DoF and  
 144 can grasp objects by the end-effector. The first 2-DoF provide the robotic arm mobility to move at the  
 145 planar 2D plane. The last DoF enables a rotational grasping by cooperating with the vision system.  
 146 The last DoF is of vital importance for the oriented object grasping, since the unrotated grasping can  
 147 easily tip the object.

148 The total takeoff weight of the UAM is about 4.08kg. Thanks to the carbon fiber material and Poly  
 149 Lactic Acid (PLA) material, the robotic arm is only weighted 0.459kg and has 43cm extended range.

### 150 2.3. Software Architecture

151 The software runs on two main processors: Pixhawk4 and Jetson TX2, and all processes are  
 152 running onboard. Fig. 5 gives an overview of the system software architecture. Note that the Jetson  
 153 TX2 is overclocked to run at the maximum clock rate and unlocked 2 external CPU cores to generate  
 154 more computing power.

155 Robot Operating System (ROS)[21] is a pseudo-operating system which allows developers to  
 156 work cooperating by following its running mechanism. Modules related to state estimation and flying

control are running on Pixhawk4 and exchange data with Jetson TX2 by MAVROS<sup>12</sup>. Visual Inertial Odometry (VIO) and object detection algorithms are running on Jetson TX2. A state machine is adapted to set the robotic arm motion and UAM waypoint. So the grasping position  $(x_w^s, y_w^s, z_w^s)$  of the UAM can be given by:

$$\begin{bmatrix} x_w^s \\ y_w^s \\ z_w^s \\ 1 \end{bmatrix} = \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix} + T_B^W T_C^B \begin{bmatrix} x_t \\ y_t \\ z_t \\ 1 \end{bmatrix} - T_B^W T_0^B \begin{bmatrix} x_0^g \\ y_0^g \\ 0 \\ 1 \end{bmatrix} \quad (1)$$

where  $T_B^W, T_0^B, T_C^B \in \mathbb{R}^{4 \times 4}$  are the homogeneous transformation matrix,  $T_0^B$  transforms the fixed arm frame to body frame,  $T_C^B$  transforms the camera frame to the body frame,  $T_B^W$  transforms the body frame to the world-fixed frame,  $(x_0^g, y_0^g)$  is the grasping point.

The standard cascaded position-velocity control of hexacopter is well implemented inside the Pixhawk4 fly controller. Reader can find more detail information at this link<sup>13</sup>. We adopt this control strategy and well tuned its inside parameters to get a robust performance.

#### 2.4. State Estimation

State estimation is the foundation of our system as it provides crucial information to help other parts to achieve the best performance. In this work, we use the Pixhawk4 built-in Extended Kalman Filter (EKF) to fuse multiple sensors feedback for state estimation.

Global Positioning System (GPS) usually fails to provide global positioning feedback when in an indoor environment. Hence, in order to fly indoor, without using expensive visual motion caption system, we integrate the VINS-mono VIO to provide the local position feedback, and it can provide the highest level of accuracy and robustness compared with multiple VIO [22]. The VINS-mono runs at 10hz with loop-closure and the output is rotated to ENU world-fixed frame denoted  $(x_w^v, y_w^v, z_w^v)$ . For the VINS-mono, we use the Inertial Measurement Unit (IMU) in Pixhawk4 as the inertial input, and a monochrome global shutter camera runs at  $640 \times 400$  resolution and 90 Frames Per Second (FPS) to provide clear images as visual input. To reduce drifting, the monochrome camera and Pixhawk4 are tight-coupled by the 3D printed structures, the camera intrinsic matrix and camera to imu transformation parameters are carefully calibrated by using kalibr[23]. Due to the computation limitation and data transmission delay, the output of VINS-mono runs in Jetson TX2 has about 140ms delay compared with current IMU output. In order to synchronize these outputs, we first calculate the velocity of VINS-mono estimation  $(\dot{x}_w^v, \dot{y}_w^v, \dot{z}_w^v)$ , then apply some random movements to the UAM, so the delay time can be found by comparing  $(\dot{x}_w^v, \dot{y}_w^v, \dot{z}_w^v)$  with the Pixhawk4 velocity estimation. Finally,  $(x_w^v, y_w^v, z_w^v)$  is used as external vision aid of the Pixhawk4 onboard EKF to give 100hz state estimation.

For robust flying, the main altitude feedback is not given by the VINS-mono but the TFmini rangefinder. The total delay of TFmini measurement is about 30ms.

<sup>12</sup> <https://github.com/mavlink/mavros>

<sup>13</sup> [https://dev.px4.io/en/flight\\_stack/controller\\_diagrams](https://dev.px4.io/en/flight_stack/controller_diagrams)

<sup>182</sup> **3. Robotic Arm & DCS**

<sup>183</sup> *3.1. Robotic Arm Motion Planning*

In order to find the workspace of robotic arm, the forward kinematics of a 3DoF robotic arm is given by the following equations:

$$\begin{aligned} x_0 &= L_1 \cos \theta_1 + L_2 \cos(\theta_1 + \theta_2) \\ y_0 &= L_1 \sin \theta_1 + L_2 \sin(\theta_1 + \theta_2) \\ \theta_3 &= \theta \end{aligned} \quad (2)$$

<sup>184</sup> where  $L_1, L_2$  are the lengths of the first and the second links.  $\theta_1, \theta_2, \theta_3$  are the rotation angles of each  
<sup>185</sup> joint,  $(x_0, y_0)$  is the point in arm fixed frame. Since the angle  $\theta_3$  is only related to the target rotation  
<sup>186</sup> angle  $\theta$ , the workspace of our robotic arm is equal to a planar 2DoF robotic arm model. The table  
<sup>187</sup> presented in Fig. 3 gives a clear D-H parameters definition of the robotic arm.

<sup>188</sup> When applying the UAM into real-world scenes, the flow generated by rotors can blow the  
<sup>189</sup> lightweight targets away, like empty plastic bottles, and it is hard to predict whether an object can  
<sup>190</sup> be easily blown away by the downward flow or not. Hence, a safety grasping action is better to be  
<sup>191</sup> taken under weak flow influenced conditions. High-fidelity Computational Fluid Dynamics (CFD)  
<sup>192</sup> simulation results of many different UAVs have been presented in [24]. From the observation of  
<sup>193</sup> these results, the flow generated by each UAV's rotor is decreasing rapidly in the outer-wing area.  
<sup>194</sup> Inspired by this observation, we plan the robotic arm to grasp in a weak flow area. Based on the flow  
<sup>195</sup> measurements described in section 5.2 and the actual mechanical movement range of robotic arm, the  
<sup>196</sup> actual workspace is given in Fig. 6, the green points indicate weak flow influenced area, red points  
<sup>197</sup> indicate strong flow influenced area. The blue star point is the dropping point and the purple star  
<sup>198</sup> point is the grasping point  $(x_0^g, y_0^g)$ . The robotic arm holds at the yellow star point  $(x_0^f, y_0^f)$  during  
<sup>199</sup> flight. And all these points can be redefined depending on the applications. We use Inverse Kinematics  
<sup>200</sup> (IK) to solve the desired rotation angle of each joint. Assumed the robotic arm is planning to move to  
<sup>201</sup>  $(x_0, y_0)$  point, the IK result is given by:

$$\theta_2 = \pm \arccos\left(\frac{x_0^2 + y_0^2 - L_1^2 - L_2^2}{2L_1L_2}\right) \quad (3)$$

<sup>202</sup> Here  $\theta_2 \in [0, \pi]$  is downward elbow solution, and  $\theta_1$  can be derived by:

$$\theta_1 = \arctan\left(\frac{x_0^2}{y_0^2}\right) - \arccos\left(\frac{x_0^2 + y_0^2 + L_1^2 - L_2^2}{2L_1\sqrt{x_0^2 + y_0^2}}\right) \quad (4)$$

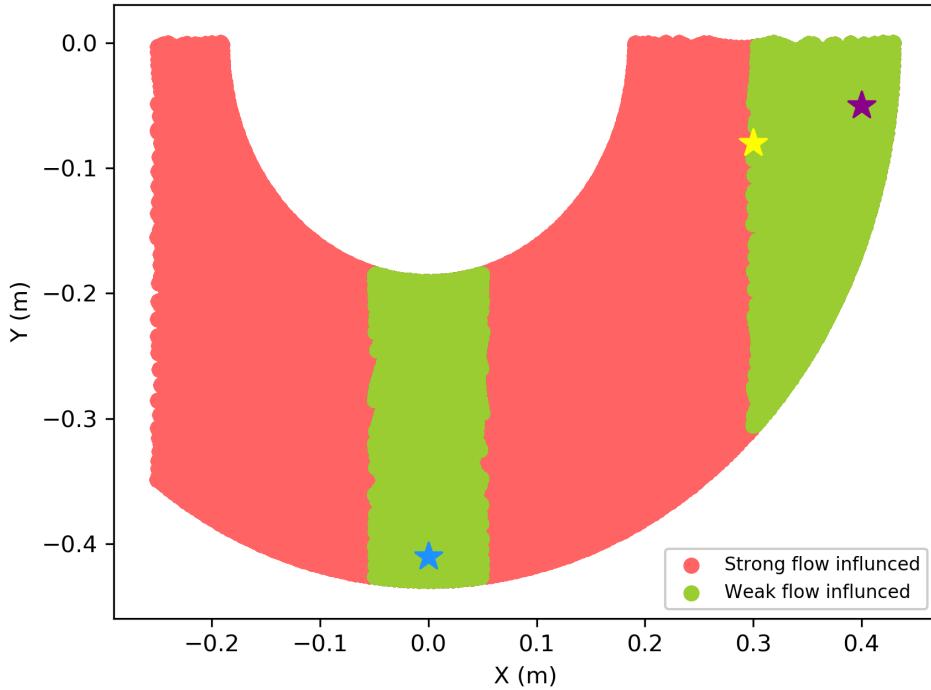
<sup>203</sup> *3.2. CoG Compensation*

<sup>204</sup> When the UAM is static on the ground and the robotic arm hold static at  $(x_0^f, y_0^f)$ , a symmetry  
<sup>205</sup> placement design is utilized to make sure the  $G_x$  and  $G_y$  are fitted with the geometry center.

<sup>206</sup> However, movements of the robotic arm can change the  $G_x$ . For the dynamic  $G_x$  alignment, we  
<sup>207</sup> adapte the strategy presented in [9] called DCS, moving the battery as a counterweight since it's weight  
<sup>208</sup> can provide sufficient compensation in the relatively short moving distance.

<sup>209</sup> CoG transformation of Link  $i$  and the end-effector payload from link frame to the body frame are  
<sup>210</sup> given by:

$$\begin{bmatrix} x_{bi}^g \\ y_{bi}^g \\ z_{bi}^g \\ 1 \end{bmatrix} = T_0^B T_i^0 \begin{bmatrix} x_i^c \\ y_i^c \\ z_i^c \\ 1 \end{bmatrix} \quad (5)$$



**Figure 6.** Workspace of robotic arm printed by using forward kinematics considered flow influence.

where  $T_i^0, T_0^B \in \mathbb{R}^{4 \times 4}$  are the homogeneous transformation matrices,  $T_0^B$  transforms from fixed arm frame to body frame,  $T_i^0$  transforms from each link frame to the fixed arm frame.  $(x_i^c, y_i^c, z_i^c)$  is the CoG position of Link  $i$  in the fixed arm frame.  $(x_{bi}^g, y_{bi}^g, z_{bi}^g)$  is the CoG position of Link  $i$  in body frame, here  $i = 3$  indicates the grasped object.

To align the  $G_x$  at geometry center, a linear slider is designed to move the battery and the position of the battery  $p_b$  in the body frame can be calculated by:

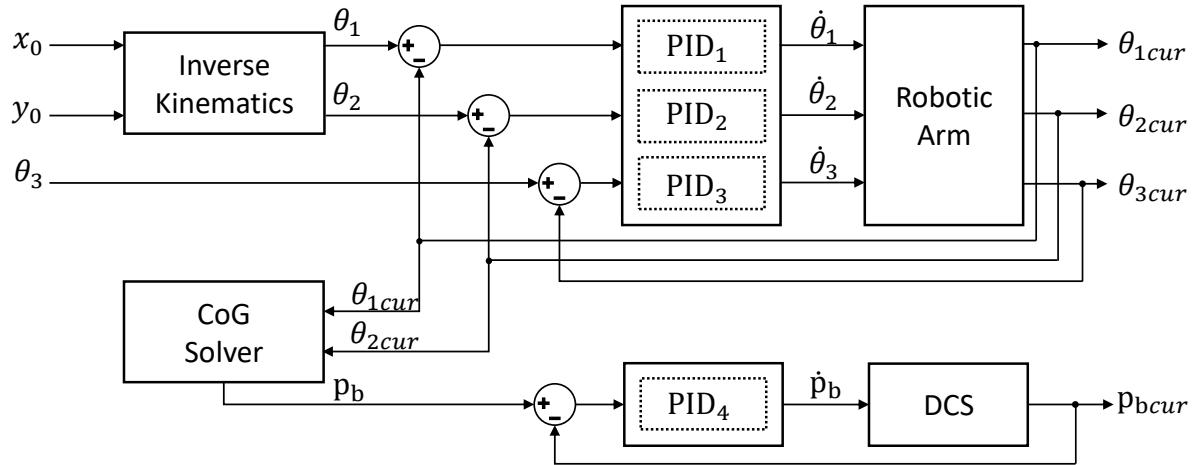
$$p_b = \frac{\sum_{i=1}^3 m_i x_{bi}^g}{m_b} \quad (6)$$

where  $m_i$  is the mass of Link  $i$ ,  $m_b$  is the mass of battery.

The displacement compensation plays a key role in stabilizing the UAM. Without the DCS, the change of CoG can easily make aside rotor reach the maximum thrust leading to unstable. And this method works well in our system. To guarantee the compensation performance, we choose a reasonable range of rotation speed of each joint in a robotic arm to make sure it does not exceed the maximum compensation speed of the linear slider.

### 3.3. Control

The total control diagram of the robotic arm and DCS are shown in Fig. 7. In Fig .7,  $\dot{\theta}_1, \dot{\theta}_2, \dot{\theta}_3, \dot{p}_b$  is the rotation speed of the servos and the  $\theta_{1cur}, \theta_{2cur}, \theta_{3cur}, p_{bcur}$  indicate the current feedback. Three PID controllers are adapted to control the robotic arm. Another PID controller uses the position of the battery  $p_b$ , detected by a laser rangefinder, to control the linear slider. All PID parameters are well tuned to guarantee stable and smooth control.



**Figure 7.** Control diagram of the DCS and robotic arm.

#### 229 4. Vision System

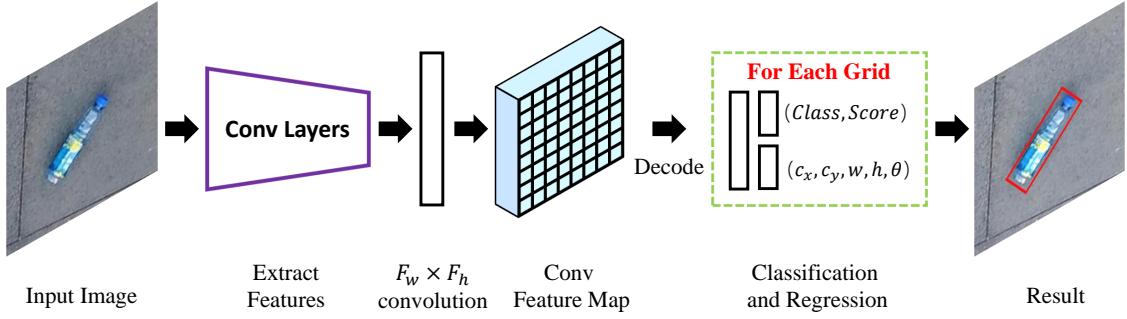
230 In this section, we introduce the vision system including the light and fast oriented-object  
 231 detection model called Rotation-SqueezeDet and the target localization framework based on  
 232 Rotation-SqueezeDet detection results and point clouds from the depth camera.

233 Our proposed model is inspired by SqueezeDet[25] and Rotation Region Proposal Networks  
 234 (RRPN)[17]. As the former, SqueezeDet is a single-pass detection pipeline combining bounding box  
 235 localization and classification by a single network. It appears to be the smallest object detector by virtue  
 236 of a powerful but small backend network of SqueezeNet[16,26]. As the latter, RRPN is based on Faster  
 237 R-CNN[14], but RRPN can detect oriented objects. The difference between Faster R-CNN and RRPN is  
 238 described as below. In Faster R-CNN, the Region of Interests (RoIs) are generated by Region Proposal  
 239 Network (RPN), and the RoIs are rectangles which can be written as  $R = (x_{min}, y_{min}, x_{max}, y_{max}) =$   
 240  $(c_x, c_y, w, h)$ . These RoIs have regressed from  $k$  anchors which are generated by some predefined  
 241 scales and aspect ratios. Then these RoIs will be feed into ROI Pooling layer and some fully connected  
 242 layers to obtain horizontal bounding boxes. However, in RRPN, instead, it uses Rotation anchors  
 243 (R-anchors) and rotation ROI Pooling, brings the ability to predict oriented bounding boxes denoted as  
 244  $R = (c_x, c_y, w, h, \theta)$ .

245 Object detection algorithms like Faster R-CNN have high accuracy but slow processing speed  
 246 and large storage requirement. Moreover, RRPN is slower than Faster R-CNN, it takes twice as much  
 247 time as the Faster-RCNN[17]. Thus RRPN does not meet our run-time requirement. Considering a  
 248 comparable accuracy and run-time on Jetson TX2, SqueezeDet is a suitable choice, it can run about  
 249 45FPS with  $424 \times 240$  pixels image on Jetson TX2 and easy to train. However, SqueezeDet cannot  
 250 predict the  $\theta$  of the rotated object because it uses horizontal bounding boxes. So we designed a model  
 251 which can generate oriented bounding boxes based on SqueezeDet and named Rotation-SqueezeDet,  
 252 it can run on Jetson TX2 in near real time.

##### 253 4.1. Network Architecture

The overall model of Rotation-SqueezeDet is illustrated in Fig. 8. In this model, a convolutional neural network, SqueezeNet V1.1, first takes an image as input and extracts a low-resolution, high dimensional feature map from the image. Then the feature map is fed into the  $F_w \times F_h$  convolutional layers to compute oriented bounding boxes at each position of conv feature map. Next, each oriented bounding box is associated with  $(5 + C + 1)$  values, where 5 is the number of bounding box parameters,  $C$  is the number of classes, and 1 is the confidence score. And each position on the convolution feature map computes  $K \times (5 + 1 + C)$  values that encode the bounding box predictions. Here,  $K$  is the



**Figure 8.** Rotation-SqueezeDet’s detection pipeline.

number of R-anchors, each R-anchor can be described by 5 parameters as  $(c_x^a, c_y^a, w^a, h^a, \theta^a)$ ,  $(c_x^a, c_y^a)$  are R-anchor’s center on image,  $w^a, h^a, \theta^a$  are the width, height and angle of R-anchor, respectively.

$$\begin{aligned} v_x &= \frac{c_x - c_x^a}{w^a}, v_y = \frac{c_y - c_y^a}{h^a}, \\ v_w &= \log \frac{w}{w^a}, v_h = \log \frac{h}{h^a}, v_\theta = \theta - \theta^a + k\pi \end{aligned} \quad (7)$$

where  $(c_x, c_y, w, h, \theta)$  are parameters describe the predicted oriented bounding box,  $(c_x^a, c_y^a, w^a, h^a, \theta^a)$  are parameters describe the R-anchor, and here  $k \in \mathbb{Z}$  to ensure  $\theta \in [0, \pi)$ .

#### 4.2. Oriented IoU Computation

To efficiently predict the rotation angle of objects, we introduce a parallel Intersection over Union (IoU) computing method. First, we attempt to calculate the IoU using the OpenCV’s functions *rotatedRectangleIntersection* and *contourArea* directly. However, the efficiency of these functions remains poor because they cannot compute parallelly. Thus, we use a simple and efficient method to approximately compute the IoU parallelly, which is to use the angle deviation of two oriented bounding boxes. The approximate IoU[27] can be computed by:

$$\text{IoU}^* = \text{IoU} * \text{abs}\left(1 - \frac{\theta_1^b - \theta_2^b}{\beta}\right) \quad (8)$$

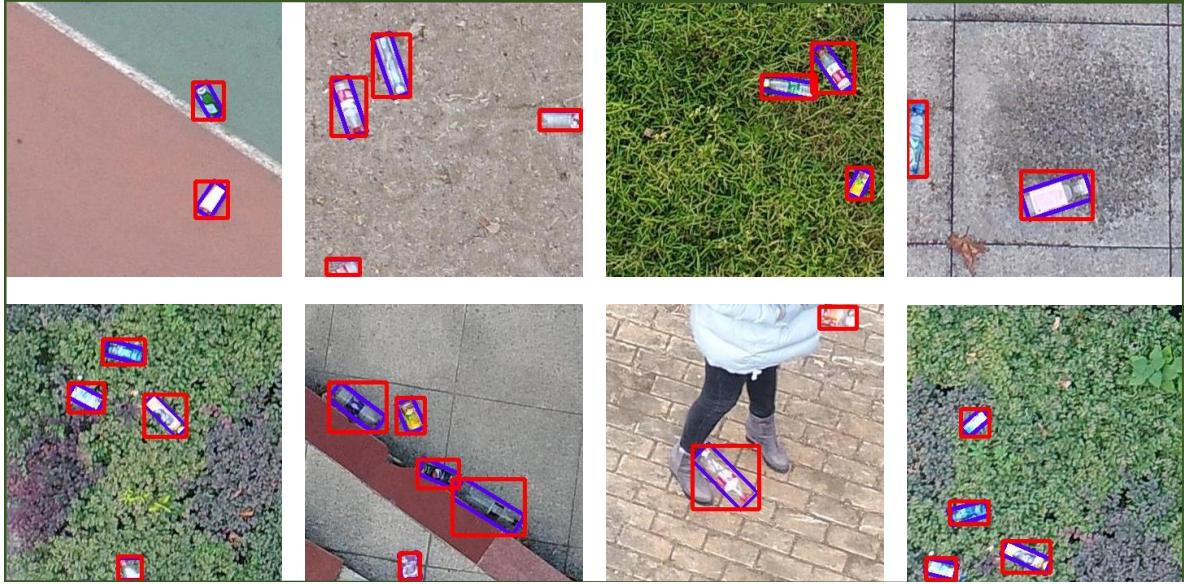
where  $\theta_1^b$  and  $\theta_2^b$  are the rotation angle of two oriented bounding boxes, IoU is computed by treating oriented bounding boxes as horizontal bounding boxes.

#### 4.3. R-Anchors’ Selection

The R-Anchors are different from the horizontal anchors. And for the R-Anchors’ selection, we use a K-means based method described in [28] to select R-anchors’  $w$  and  $h$  to match the data distribution, we set  $k$  as 9 in K-means and treat objects’ angle distribution as a uniform distribution, i.e., we set R-anchors’ angle as  $\{0, \frac{\pi}{9}, \frac{2\pi}{9}, \frac{3\pi}{9}, \frac{4\pi}{9}, \frac{5\pi}{9}, \frac{6\pi}{9}, \frac{7\pi}{9}, \frac{8\pi}{9}\}$ . Therefore there are 81 anchors at each conv feature map position.

#### 4.4. Object Localization

To acquire the real world position of the target, we use the RGB-D camera to detect and locate the target. First, the aligned RGB image and point clouds can be obtained by aligning RGB image and depth image in the same coordinate system. Next, the subarea of the total point clouds containing location information of the target which can be extracted from the whole point clouds by utilizing the



**Figure 9.** Detection results of the proposed detection algorithm in UAV-BD dataset. Our algorithm can give the common bounding box in red and oriented bounding box in blue at the same time.

276 detection result  $(c_x, c_y, w, h, \theta)$ . After that, we use a small central subarea of target's point clouds to  
277 calculate its real-world position  $(x_t, y_t, z_t)$  in the camera frame given by:

$$(x_t, y_t, z_t) = \left( \frac{1}{L} \sum_{i=0}^{k^2} \mathbf{X}_p^i, \frac{1}{M} \sum_{i=0}^{k^2} \mathbf{Y}_p^i, \frac{1}{N} \sum_{i=0}^{k^2} \mathbf{Z}_p^i \right) \quad (9)$$

where  $(\mathbf{X}_p, \mathbf{Y}_p, \mathbf{Z}_p)$  are the coordinates of the point clouds set of the center subarea of target's bounding box, its superscript indicates the  $i^{th}$  point value started from top left corner.  $L, M, N$  are the valid points number of  $\mathbf{X}_p, \mathbf{Y}_p, \mathbf{Z}_p$ , respectively. The central subarea of target's point clouds can be written as  $(c_x, c_y, k, k)$ , the  $k \times k$  is the size of the selected central subarea of target's point clouds, which can be calculated by:

$$k = \begin{cases} 5 & \text{if } \min(w, h) > 5 \\ \min(w, h) & \text{otherwise} \end{cases}$$

278 Here we set  $k$  as 5 to reduce the computing burden. Finally, the parameters of position and rotation  
279 angle of a target can be written as:  $(x_t, y_t, z_t, \theta)$ ,  $\theta$  is the rotation angle of target's anchor.

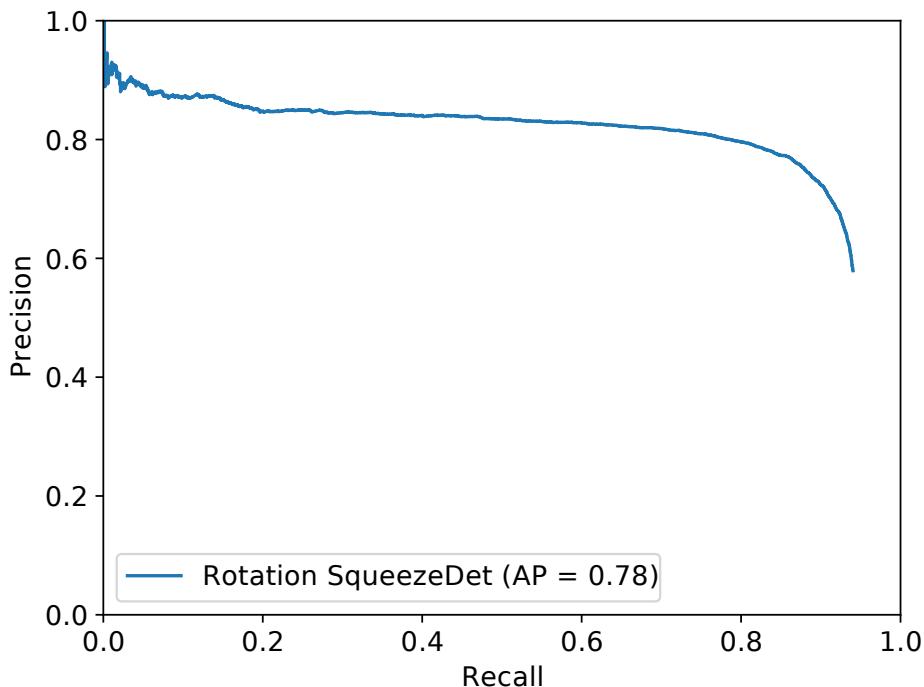
## 280 5. Experiments

281 The experimental setup is mentioned in Section 2. During the flight tests, all data are logged  
282 onboard with no external data transmission.

### 283 5.1. Vision System Results

284 We trained and evaluated our vision system based on our previous work UAV-BD[29], a bottle  
285 image dataset under UAV perspective. It contains about 34,791 object instances in 25,407 images  
286 labeled by oriented bounding boxes. For training and evaluating our model, 64% of the images were  
287 randomly selected as the training data, 16% as validation data and the rest 20% as the testing data.

288 All object detection experiments were implemented on TensorFlow[30]. We used the pertained  
289 model, SqueezeNet v1.1, to initialize the network. And the system was trained 100k steps with a batch  
290 size of 20 and a learning rate of 0.01. Besides, weight decay and momentum were 0.0001 and 0.9,  
291 respectively. The optimizer was *MomentumOptimizer*.



**Figure 10.** AP results of Rotation-SqueezeNet validated in UAV-BD.

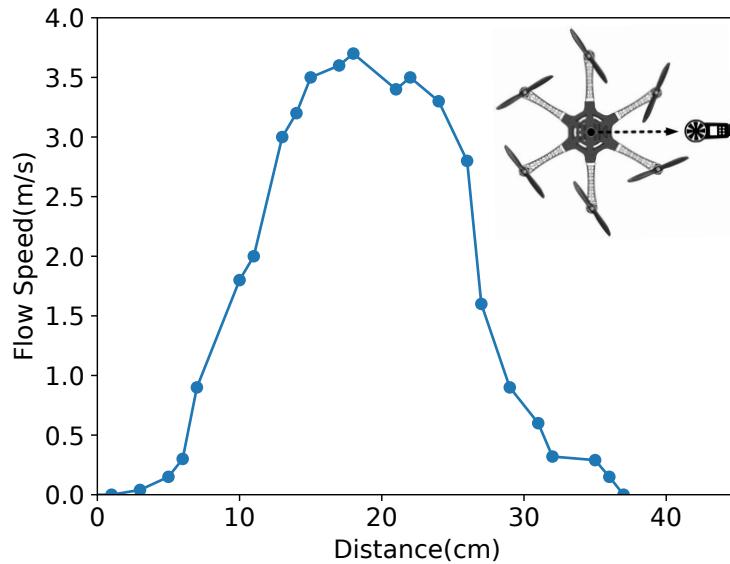
As shown in Fig. 10, the Average Precision (AP) of Rotation-SqueezeDet on UAV-BD is about 78.0% when IoU = 0.5. The run-time on Jetson TX2 is about 41ms with the size of  $424 \times 240$  pixels color image. In Fig. 9 we visualize the detection results in UAV-BD dataset. Our algorithms can be given two results, first is the common horizontal bounding box, marked by red color, without any rotation, second is the oriented blue bounding box inside the red bounding box. We can clearly see that the blue boxes are more accurate than the red boxes since it contains fewer background pixels and locates the object more precisely.

### 5.2. Flow Distribution Validation

In order to give a rough parameter estimation for the planning mentioned in Section 3, we disarmed the UAM and keep it flying at a certain height about 1.5 m above the ground. Then use Smart AS856<sup>14</sup> anemometer to measure the downward flow distribution. Since the flow is highly complex, we cannot give a very accurate flow distribution by just using an anemometer, thus we only measured a 1D flow distribution to give a rough estimation for the robotic arm planning. The moving path of the anemometer is given in the top right corner of Fig. 11 and the vertical distance from the anemometer to the rotor is about 16cm to represent the overall rough distribution below the UAM.

Following this path, we recorded the average value of 5 measurements at one point and drew the curve in Fig. 11. The horizontal axis is the distance from the central position of the UAM body to the outside following this path. From the observation of Fig. 11, the flow speed is decreasing rapidly at about 8cm and 30cm, and reaching the top speed at about 21cm which is the underside of rotors. The 8cm is close to the inside of the UAM, the 30cm is close to the outside of the UAM. So the flow speed is relatively weak when at the central position and outside position of the UAM.

<sup>14</sup> [en.smartsensor.cn/products\\_detail/productId=248.html](http://en.smartsensor.cn/products_detail/productId=248.html)



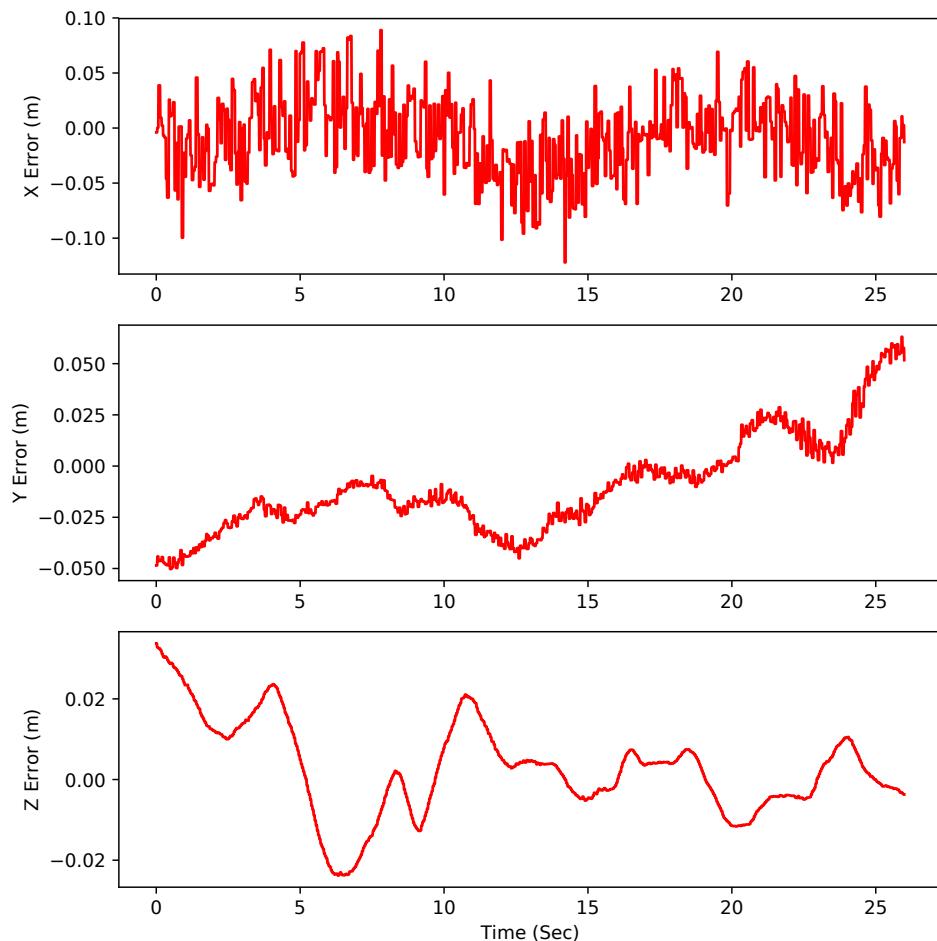
**Figure 11.** The flow speed measured by digital anemometer and the subfigure in the top right corner is the path of the anemometer during measurement.

### 313 5.3. Aerial Robotic Arm Moving Test

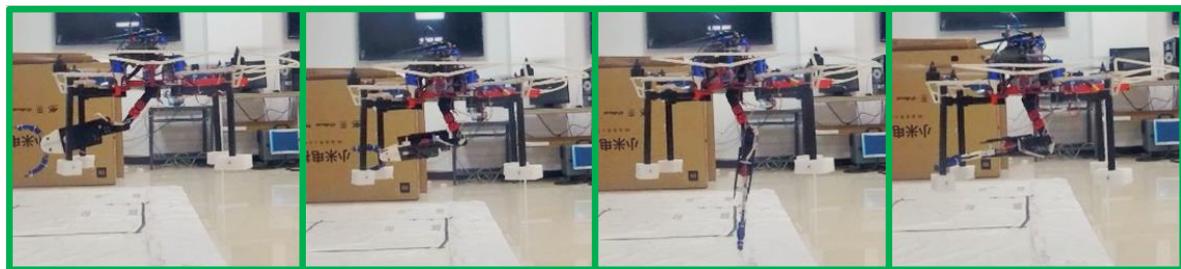
314 To evaluate the stability and controllability of our system, the UAV was programmed to hovering  
 315 at a certain point. In Fig. 12, the control errors of the system during this flight were logged and  
 316 presented. After finishing the takeoff procedure, the robotic arm will automatically move to multiple  
 317 points and rotate the end-effector  $90^\circ$ . The DCS was activated in this flight, moving the battery to  
 318 compensate for the change of CoG generated by movements of the robotic arm. The standard deviation  
 319 of control error in the world-fixed frame is about 3.64cm in x axis, 2.37cm in y axis and 1.16cm in z  
 320 axis, indicated our system can keep hovering at the certain point no matter the robotic arm is moving  
 321 or not. We noticed that during the whole test, the error increased at several points including 5, 10 and  
 322 22 seconds, respectively. This was caused by the moving of the robotic arm, but our system could  
 323 always stabilize itself. Fig. 13 shows some snapshots of the robotic arm aerial moving experiments.

### 324 5.4. Autonomous Grasping

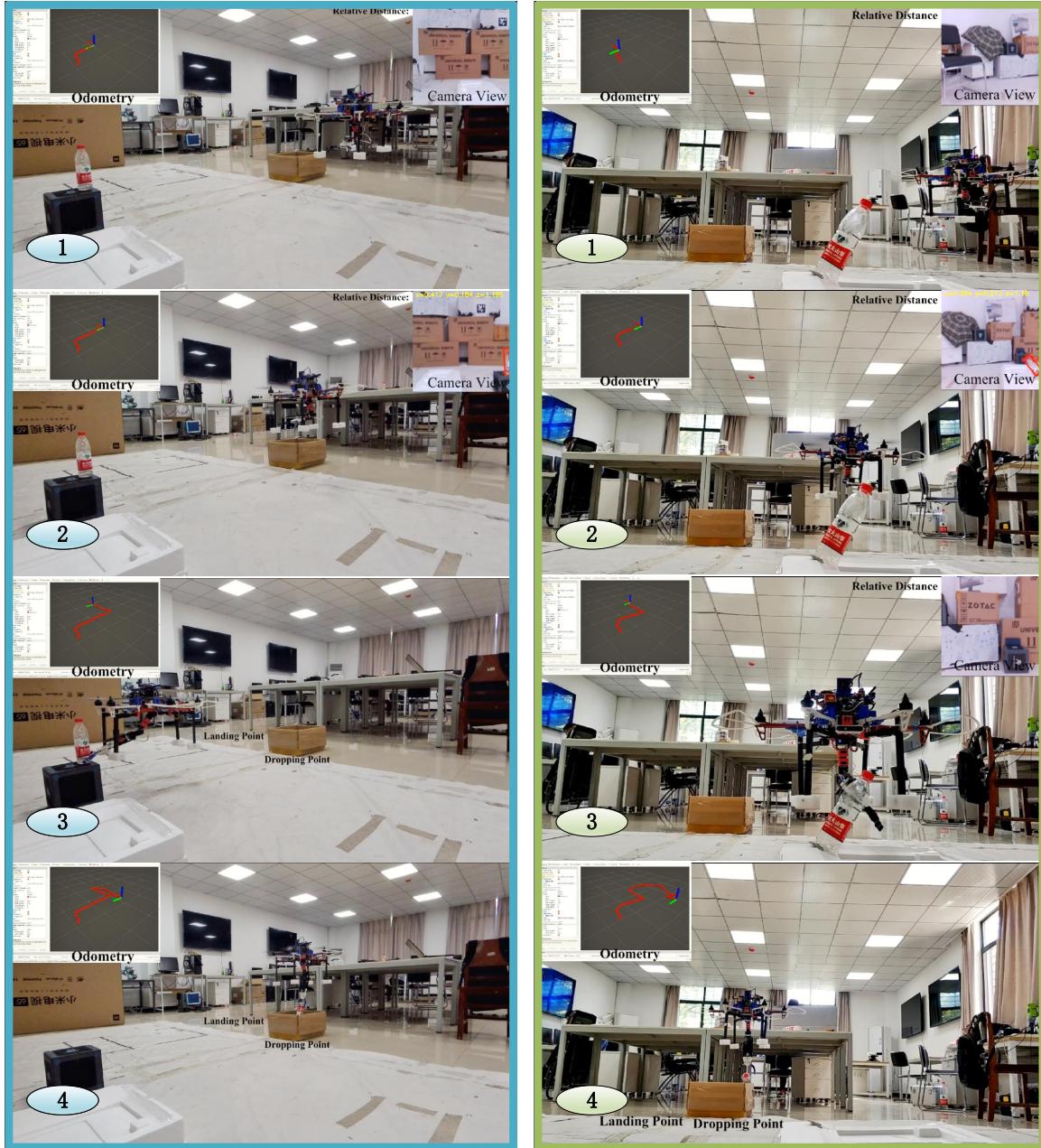
325 Autonomous grasping experiments with objects placed vertically and obliquely have also been  
 326 conducted. Fig. 14 presents the key steps of these two grasping experiments. The subfigures were  
 327 taken at the specific moments ordered by the number in the bottom left corner: ① indicated the UAM  
 328 is searching specific target. Here we use the plastic bottle as the target; ② indicated the UAM has  
 329 detected the target and given its relative position. The UAM will then align its position to the grasping  
 330 position; ③ indicated the UAM hovering at the grasping position and the end effector is about to grasp  
 331 the bottle. ④ indicated the UAM dropping the bottle at the dropping point. The odometry given by  
 332 the onboard EKF in 100hz rates were simultaneously presented in the top left corner, showing our  
 333 system can run indoor without using any external visual motion caption system. The detection results  
 334 and the calculated relative distance of the target objects were also simultaneously presented in the top  
 335 right corner during searching. In Fig. 14(a), the target was vertically placed at the top of a cube. By  
 336 comparing images captured at ① and ②, it is easily to know the bottle is not shown in ① but detected in  
 337 ②. The UAM was using a constant speeded to search in ①, and speeded up to get close to target after  
 338 ②. The end-effector successfully approached and grabbed the target in ③ verified our methods worked  
 339 well, and the plastic bottle was not flipped by the downward flow. In ④ the UAM can automatically



**Figure 12.** Control errors history during the whole aerial robotic arm moving.



**Figure 13.** Snapshots of the robotic arm aerial moving experiments. Many movements are carried out during this flight to test the system.



**Figure 14.** Snapshots of autonomous grasping experiments. The number in subfigures indicate selected moments: ① indicated the UAM is searching specific target; ② indicated the UAM detected the target and gives its relative position; ③ indicated the UAM hovering at the grasping position and the end effector is about to grasp the bottle. ④ indicated the UAM dropping the grasped target at the dropping points.

340 drop the target at the dropping point indicates our system is capable of automatically finishing the  
341 whole grasping task.

342 In Fig. 14(b), the procedures were basically the same as mentioned above. What different is that  
343 we placed the empty bottle with about 45° rotation. ② gave the detected target with its relative position  
344 and rotated angle. In ③ the end-effector used the rotation information of the target to successfully  
345 grasp the target.

346 These experimental results have shown that our system can accomplish autonomous an grasping  
347 mission.

## 348 6. Conclusion

349 In this paper, we developed an approach to enable a UAM to grasp oriented objects. The key  
350 challenges included detecting and locating objects in UAM perspective, CoG compensation, the flow  
351 influence to the objects and indoor positioning of UAM system, all of which made the autonomous  
352 grasping mission very difficult. We showed the effectiveness of our approach in real tests with the  
353 ability to detect, locate and grasp objects with arbitrary poses in GPS-denied environments without  
354 relying on the visual motion caption system. We believe that the proposed solution, both in terms  
355 of hardware and algorithms, will be useful in not only the aerial grasping missions but also general  
356 grasping missions since the vision system can be applied to any kind of robotic arm. Future work will  
357 be set out to investigate how to estimate the 6D pose of objects in real-time. We will also improve the  
358 stability by using the servos with torque feedback and adopting a more advanced controller for the  
359 UAM.

360 **Author Contributions:** Conceptualization, S.L. and J.W.; methodology, S.L. and J.W.; software, S.L., J.W. and R.P.;  
361 validation, S.L., J.W. and R.P.; writing—original draft preparation, S.L. and J.W.; writing—review and editing, W.Y.;  
362 supervision, W.Y.; project administration, S.L. and W.Y.

363 **Funding:** The work was supported in part by the Funds for Creative Research Groups of Natural Science  
364 Foundation of Hubei Province under Grant 2018CFA006.

365 **Conflicts of Interest:** The authors declare no conflict of interest.

## 366 References

- 367 1. Lendzioch, T.; Langhammer, J.; Jenicek, M. Estimating Snow Depth and Leaf Area Index Based on UAV  
368 Digital Photogrammetry. *Sensors* **2019**, *19*. doi:10.3390/s19051027.
- 369 2. Olivares-Mendez, M.A.; Fu, C.; Ludivig, P.; Bissyandé, T.F.; Kannan, S.; Zurad, M.; Annaiyan, A.; Voos, H.;  
370 Campoy, P. Towards an Autonomous Vision-Based Unmanned Aerial System against Wildlife Poachers.  
*Sensors* **2015**, *15*, 31362–31391. doi:10.3390/s151229861.
- 372 3. Ruggiero, F.; Lippiello, V.; Ollero, A. Aerial Manipulation: A Literature Review. *IEEE Robotics and Automation  
373 Letters* **2018**, *3*, 1957–1964. doi:10.1109/LRA.2018.2808541.
- 374 4. Kim, S.; Choi, S.; Kim, H.J. Aerial manipulation using a quadrotor with a two DOF robotic arm. 2013  
375 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2013, pp. 4990–4995.  
376 doi:10.1109/IROS.2013.6697077.
- 377 5. Korpela, C.; Orsag, M.; Danko, T.; Oh, P. Insertion tasks using an aerial manipulator. 2014 IEEE  
378 International Conference on Technologies for Practical Robot Applications (TePRA), 2014, pp. 1–6.  
379 doi:10.1109/TePRA.2014.6869148.
- 380 6. Božek, P.; Al Akkad M, A.; Blištan, P.; Ibrahim N, I. Navigation control and stability investigation of a mobile  
381 robot based on a hexacopter equipped with an integrated manipulator. *International Journal of Advanced  
382 Robotic Systems* **2017**, *14*, 1–13.
- 383 7. Arleo, G.; Caccavale, F.; Muscio, G.; Pierri, F. Control of quadrotor aerial vehicles equipped with a robotic  
384 arm. 21St Mediterranean Conference on Control and Automation (MED). IEEE, 2013, pp. 1174–1180.
- 385 8. Jimenez-Cano, A.; Martin, J.; Heredia, G.; Ollero, A.; Cano, R. Control of an aerial robot with multi-link arm  
386 for assembly tasks. 2013 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2013, pp.  
387 4916–4921.

- 388 9. Ruggiero, F.; Trujillo, M.A.; Cano, R.; Ascorbe, H.; Viguria, A.; Peréz, C.; Lippiello, V.; Ollero, A.; Siciliano,  
389 B. A multilayer control for multirotor UAVs equipped with a servo robot arm. 2015 IEEE International  
390 Conference on Robotics and Automation (ICRA). IEEE, 2015, pp. 4014–4020.
- 391 10. Ohnishi, Y.; Takaki, T.; Aoyama, T.; Ishii, I. Development of a 4-joint 3-DOF robotic arm with anti-reaction  
392 force mechanism for a multicopter. 2017 IEEE/RSJ International Conference on Intelligent Robots and  
393 Systems (IROS), 2017, pp. 985–991. doi:10.1109/IROS.2017.8202265.
- 394 11. Kim, S.; Seo, H.; Choi, S.; Kim, H.J. Vision-Guided Aerial Manipulation Using a Multirotor With a Robotic  
395 Arm. *IEEE/ASME Transactions on Mechatronics* **2016**, *21*, 1912–1923. doi:10.1109/TMECH.2016.2523602.
- 396 12. Ramon Soria, P.; Arrue, B.C.; Ollero, A. Detection, location and grasping objects using a stereo sensor on  
397 UAV in outdoor environments. *Sensors* **2017**, *17*, 103.
- 398 13. Kanellakis, C.; Terreran, M.; Kominik, D.; Nikolakopoulos, G. On vision enabled aerial manipulation for  
399 multirotors. 2017 22nd IEEE International Conference on Emerging Technologies and Factory Automation  
400 (ETFA), 2017, pp. 1–7. doi:10.1109/ETFA.2017.8247653.
- 401 14. Ren, S.; He, K.; Girshick, R.; Sun, J. Faster r-cnn: Towards real-time object detection with region proposal  
402 networks. Advances in neural information processing systems, 2015, pp. 91–99.
- 403 15. Redmon, J.; Farhadi, A. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767* **2018**.
- 404 16. Iandola, F.N.; Han, S.; Moskewicz, M.W.; Ashraf, K.; Dally, W.J.; Keutzer, K. SqueezeNet: AlexNet-level  
405 accuracy with 50x fewer parameters and <0.5MB model size. *arXiv:1602.07360* **2016**.
- 406 17. Ma, J.; Shao, W.; Ye, H.; Wang, L.; Wang, H.; Zheng, Y.; Xue, X. Arbitrary-Oriented Scene Text Detection via  
407 Rotation Proposals. *IEEE Transactions on Multimedia* **2018**, *20*, 3111–3122. doi:10.1109/TMM.2018.2818020.
- 408 18. Yang, X.; Sun, H.; Fu, K.; Yang, J.; Sun, X.; Yan, M.; Guo, Z. Automatic ship detection in remote sensing  
409 images from google earth of complex scenes based on multiscale rotation dense feature pyramid networks.  
*Remote Sensing* **2018**, *10*, 132.
- 410 19. Qin, T.; Li, P.; Shen, S. Vins-mono: A robust and versatile monocular visual-inertial state estimator. *IEEE  
Transactions on Robotics* **2018**, *34*, 1004–1020.
- 411 20. Siciliano, B.; Sciavicco, L.; Villani, L.; Oriolo, G. *Robotics: modelling, planning and control*; Springer Science &  
412 Business Media, 2010.
- 413 21. Quigley, M.; Conley, K.; Gerkey, B.; Faust, J.; Foote, T.; Leibs, J.; Wheeler, R.; Ng, A.Y. ROS: an open-source  
414 Robot Operating System. ICRA workshop on open source software. Kobe, Japan, 2009, Vol. 3, p. 5.
- 415 22. Delmerico, J.; Scaramuzza, D. A Benchmark Comparison of Monocular Visual-Inertial Odometry Algorithms  
416 for Flying Robots. *Memory* **2018**, *10*, 20.
- 417 23. Rehder, J.; Nikolic, J.; Schneider, T.; Hinzmann, T.; Siegwart, R. Extending kalibr: Calibrating the extrinsics  
418 of multiple IMUs and of individual axes. 2016 IEEE International Conference on Robotics and Automation  
419 (ICRA). IEEE, 2016, pp. 4304–4311.
- 420 24. Diaz, P.V.; Yoony, S. High-Fidelity Computational Aerodynamics of Multi-Rotor Unmanned Aerial Vehicles.  
421 2018 AIAA Aerospace Sciences Meeting, 2018, p. 1266.
- 422 25. Wu, B.; Wan, A.; Iandola, F.; Jin, P.H.; Keutzer, K. SqueezeDet: Unified, Small, Low Power Fully  
423 Convolutional Neural Networks for Real-Time Object Detection for Autonomous Driving. 2017 IEEE  
424 Conference on Computer Vision and Pattern Recognition (CVPR), 2017.
- 425 26. Tripathi, S.; Dane, G.; Kang, B.; Bhaskaran, V.; Nguyen, T. LCDet: Low-Complexity Fully-Convolutional  
426 Neural Networks for Object Detection in Embedded Systems. 2017 IEEE Conference on Computer Vision  
427 and Pattern Recognition Workshops (CVPRW), 2017, pp. 411–420. doi:10.1109/CVPRW.2017.56.
- 428 27. Xie, L.; Ahmad, T.; Jin, L.; Liu, Y.; Zhang, S. A New CNN-Based Method for Multi-Directional Car License  
429 Plate Detection. *IEEE Transactions on Intelligent Transportation Systems* **2018**, *19*, 507–517.
- 430 28. Redmon, J.; Farhadi, A. YOLO9000: Better, Faster, Stronger. 2017 IEEE Conference on Computer Vision and  
431 Pattern Recognition (CVPR) **2017**, pp. 6517–6525.
- 432 29. Wang, J.; Guo, W.; Pan, T.; Yu, H.; Duan, L.; Yang, W. Bottle Detection in the Wild Using Low-Altitude  
433 Unmanned Aerial Vehicles. 2018 21st International Conference on Information Fusion (FUSION), 2018, pp.  
434 439–444. doi:10.23919/ICIF.2018.8455565.
- 435 30. Abadi, M.; Barham, P.; Chen, J.; Chen, Z.; Davis, A.; Dean, J.; Devin, M.; Ghemawat, S.; Irving, G.; Isard, M.;  
436 others. Tensorflow: a system for large-scale machine learning. 12th USENIX Symposium on Operating  
437 Systems Design and Implementation (OSDI), 2016, Vol. 16, pp. 265–283.
- 438 439

**440** © 2019 by the authors. Submitted to *Sensors* for possible open access publication under the terms and conditions  
**441** of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).