

PRGFlow: Benchmarking SWAP-Aware Unified Deep Visual Inertial Odometry

Nitin J. Sanket, Chahat Deep Singh, Cornelia Fermüller, Yiannis Aloimonos

Abstract—Odometry on aerial robots has to be of low latency and high robustness whilst also respecting the Size, Weight, Area and Power (SWAP) constraints as demanded by the size of the robot. A combination of visual sensors coupled with Inertial Measurement Units (IMUs) has proven to be the best combination to obtain robust and low latency odometry on resource-constrained aerial robots. Recently, deep learning approaches for Visual Inertial fusion have gained momentum due to their high accuracy and robustness. However, the remarkable advantages of these techniques are their inherent scalability (adaptation to different sized aerial robots) and unification (same method works on different sized aerial robots) by utilizing compression methods and hardware acceleration, which have been lacking from previous approaches.

To this end, we present a deep learning approach for visual translation estimation and loosely fuse it with an Inertial sensor for full 6DoF odometry estimation. We also present a detailed benchmark comparing different architectures, loss functions and compression methods to enable scalability. We evaluate our network on the MSCOCO dataset and evaluate the VI fusion on multiple real-flight trajectories.

Keywords – *Deep Learning in Robotics and Automation, Aerial Systems: Perception and Autonomy, Sensor Fusion, SLAM.*

I. INTRODUCTION

A fundamental competence of aerial robots [1] is to estimate ego-motion or odometry before any control strategy is employed [2]–[5]. Different sensor combinations have been used previously to aid the odometry estimation with LIDAR based approaches topping the accuracy charts [6], [7]. However such approaches cannot be used on smaller aerial robots due to their size, weight and power constraints. Such small aerial robots are generally preferred due to safety, agility and usability as swarms [8]–[10]. In the last decade, imaging sensors have struck the right balance considering accuracy and general sensor utility [11]. However, visual data is dense and requires a lot of computation, which creates challenges for low-latency applications. To this end, sensor fusion experts proposed to use IMUs along with imaging sensors, because IMUs are lightweight and are generally available on aerial robots [12]. Also, employing IMUs with even a monocular camera enables the estimation of metric depth which can be useful for many applications.

In the last decade, several VIO algorithms have been used in commercial products and also many algorithms have been made open-source by the research community. However, there is no trivial way of downscaling these algorithms for smaller aerial robots [13].

In the last five years, deep learning based approaches for visual and visual inertial odometry estimation have gained

momentum. We classify as such algorithms all approaches which learn to predict odometry in an end-to-end fashion using one of the aforementioned sensors or which use deep learning as a part of the odometry estimation. The networks used in these approaches can be compressed to smaller size with generally a linear drop in accuracy to cater to SWAP constraints. The critical issue with deep networks for odometry estimation is that to have the same accuracy as classical approaches they are generally computationally heavy leading to larger latencies. However, leveraging hardware acceleration and better parallelizable architectures can mitigate this problem.

In this work, we present a method for visual inertial odometry estimation targeted towards a down-facing/up-facing camera coupled to an altimeter source such as a barometer (outdoor) or SONAR or single beam LIDAR (indoor). Our approach uses deep learning to obtain translation - shift and zoom-in/out and/or yaw. The inputs to the network are rotation compensated using Inertial estimates of attitude. Finally, we use the altimeter to scale the shifts to real-world velocities similar to [14]. We also benchmark different combinations of our approach and answer the following questions: *How many warping blocks to use? What network architecture to use? Which loss function to use? What is the best way to compress? Which common hardware is the best for a certain-sized aerial robot?*

A. Problem Definition and Contributions

A quadrotor is equipped with a down/up-facing camera coupled to an altimeter and an IMU. The aim is to estimate ego-motion or odometry combining all sources of information. The presented approach has to be *scalable* and *unified* so that the same method can be used on different sized aerial robots catering to different SWAP constraints (Fig. 1 shows examples of different sized quadrotors with different components which can be used on quadrotors).

A summary of our contributions is given below:

- A deep learning approach to estimate odometry using visual, inertial and altimeter data
- A comprehensive benchmark of different network architectures, hardware architectures and loss functions
- Real-flight experiments demonstrating robustness of the presented approach
- Notes to practitioners whenever applicable

B. Related Work

There has been extensive progress in the field of Visual or Visual-Inertial Odometry (VI or VIO) using classical

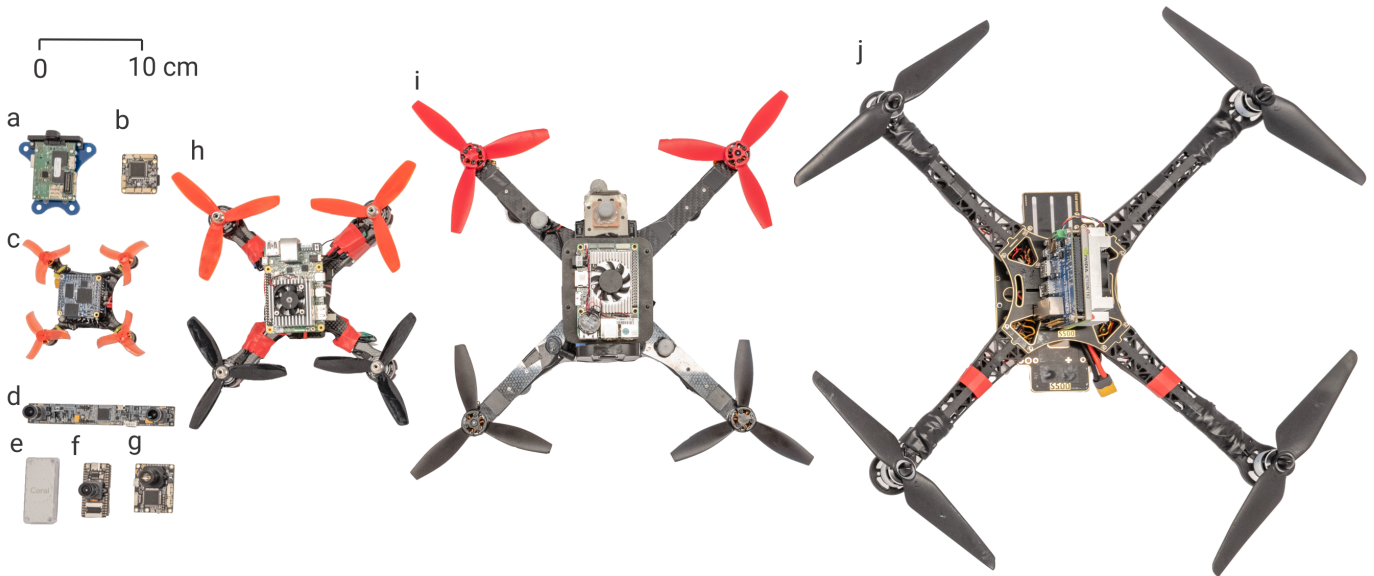


Figure 1. Size comparison of various components used on quadrotors. (a) Snapdragon Flight, (b) PixFalcon, (c) 120 mm quadrotor platform with NanoPi Neo Core 2, (d) MYNT EYE stereo camera, (e) Google Coral USB accelerator, (f) Sipeed Maix Bit, (g) PX4Flow, (h) 210 mm quadrotor platform with Coral Dev board, (i) 360 mm quadrotor platform with Intel[®] Up board, (j) 500 mm quadrotor platform with NVIDIA[®] Jetson[™] TX2. Note that all components shown are to relative scale. *All the images in this paper are best viewed in color.*

approaches, but adapting them to deep learning is still in a nascent stage. We categorize related work into the following three categories: VI/VIO using classical methods (non-deep learning), deep learning based VI/VIO, and deep learning and odometry benchmarks. Also, note that we do not consider Simultaneous Localization And Mapping (SLAM) approaches [15] such as ORB-SLAM [16], LSD-SLAM [17], LOAM [6], V-LOAM [7] and probabilistic object centric slam [18]. We also exclude LIDAR and SONAR based odometry approaches from our discussion.

1) *VI/VIO using classical methods*: Following are the state-of-the-art approaches in chronological order.

- MSCKF [12] proposed an Extended Kalman Filter (EKF) for visual inertial odometry.
- OKVIS [19] proposed a stereo-keyframe based sliding window estimator to reduce landmark re-projection errors.
- ROVIO [20] also uses an EKF but included tracking 3D landmarks along with tracking of image patches.
- DSO [21] uses a direct approach using a photometric model coupled with a geometric model to achieve the best compromise of speed and accuracy.
- VINS-Mono [22] introduced a non-linear optimization based sliding window estimator with pre-integrated IMU factors.
- Saliency-DSO [23] builds upon DSO to add visual saliency using deep learning for feature extraction. However, the optimization or regression of poses is performed classically.

2) *Deep learning based VI/VIO*:

- PoseNet [24] uses a deep network to re-localize a camera in a pre-trained scene which brought the robustness and

ease of use of deep networks for camera pose regression into limelight. Better loss functions for the same function were presented in [25].

- SfMLearner [26] took this one step further to regress camera poses and depth simultaneously from a sequence of video frames in a completely self-supervised (unsupervised) manner using geometric constraints.
- GeoNet [27] built upon SfMLearner to add additional geometric constraints and proposed a new training method along with a novel network architecture to predict pose, depth and optical flow in a completely self-supervised (unsupervised) manner.
- D3VO [28] tightly incorporates the predicted depth, pose and uncertainty into a direct visual odometry method to boost both the front-end tracking as well as the back-end non-linear optimization.
- VINet [29] proposed a supervised method to estimate odometry from a CNN + LSTM combination using both visual and inertial data. This approach, however, does not present results about its generalizability to novel scenes.
- DeepVIO [30] presents an approach to tightly integrate visual and inertial features using a CNN + LSTM to estimate odometry. This method also does not present results about its generalizability to novel scenes.

3) *Deep learning and Odometry benchmarks*: A myriad of datasets such as KITTI [31], EuRoC [32], TUMMonoVO [33], and PennCOSYVIO [34] have been proposed to evaluate the performance of VI/VIO approaches, but they do not contain enough images to train a neural network to generalize to novel datasets.

Though there exist several benchmarks for either classical VO/VIO approaches [35] and for deep learning for

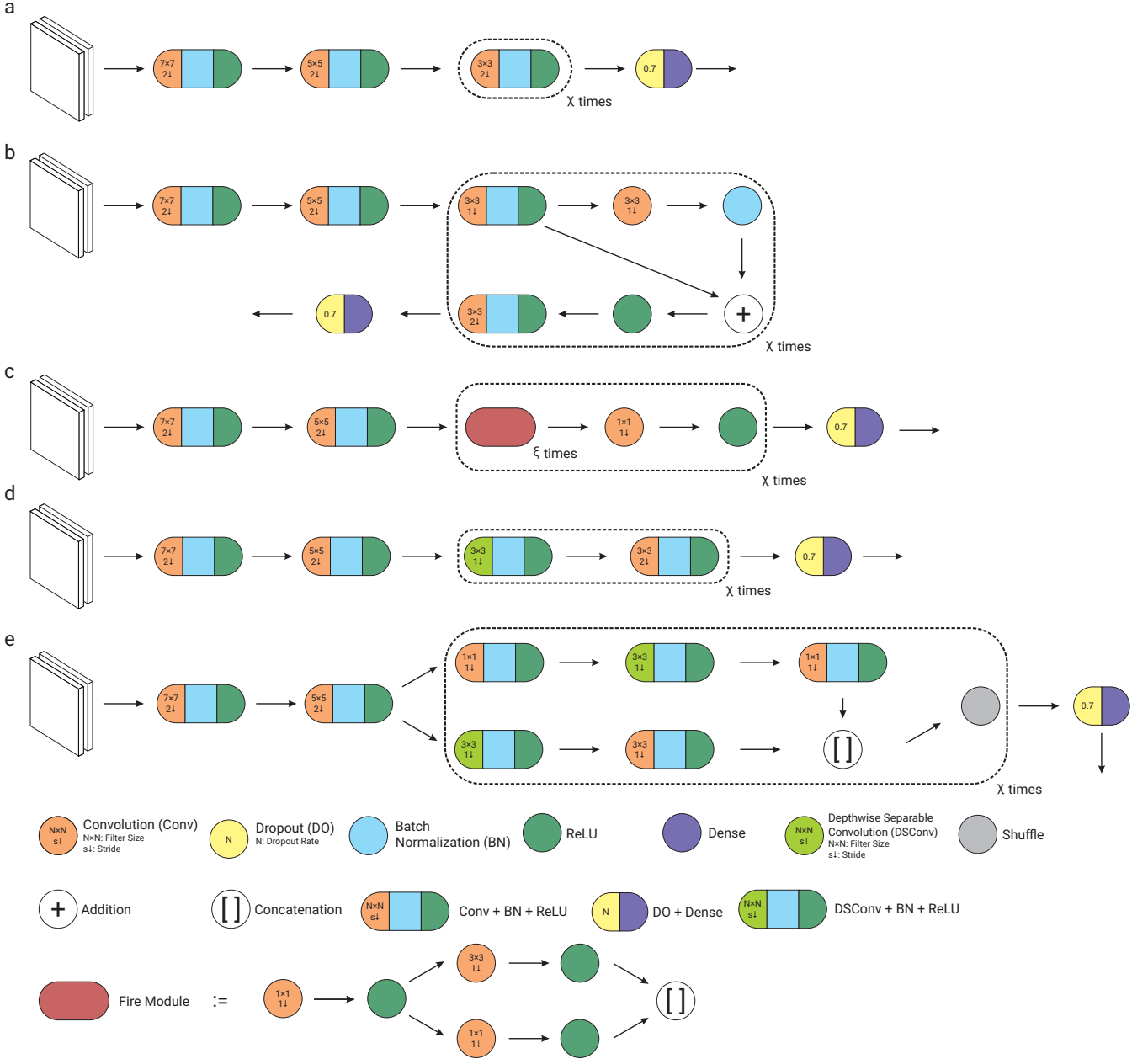


Figure 2. Different network architectures. (a) VanillaNet, (b) ResNet, (c) SqueezeNet, (d) MobileNet and (e) ShuffleNet. (χ and ξ are hyperparameters). Each architecture block is repeated per warp parameter prediction. *This image is best viewed on the computer screen at a zoom of 200%.*

classification/regression tasks [36]–[38], there is a big void in benchmarks for deep learning based VO/VIO approaches, which is the focus of this work.

II. PSEUDO-SIMILARITY ESTIMATION USING PRGFLOW

Let us mathematically formulate our problem statement. Let \mathcal{I}_t and \mathcal{I}_{t+1} be the image frames captured at times t and $t + 1$, respectively. Now, the transformation between the image frames can be expressed as $\mathbf{x}_{t+1} = \mathbf{H}_t^{t+1} \mathbf{x}_t$, where $\mathbf{x}_{t+1}, \mathbf{x}_t$ represent the homogeneous point correspondences in the two image frames and \mathbf{H}_t^{t+1} is the non-singular 3×3 transformation matrix between the two frames.

While in general the transformation between views is a non-linear function of the 3D rotation matrix \mathbf{R}_t^{t+1} and 3D translation vector \mathbf{T}_t^{t+1} , for certain scene structures, the transformation can simplify to a linear function. One such case is when the real world area is planar or can be approximated to be a plane, or the focal length is large. This scenario is also called a homography with \mathbf{H}_t^{t+1} referred to as the *Homography matrix*.

From \mathbf{H}_t^{t+1} we can recover a finite number of $\{\mathbf{R}_t^{t+1}, \mathbf{T}_t^{t+1}\}$ solutions. Furthermore, \mathbf{H}_t^{t+1} may also be decomposed into simpler transformations such as in-plane rotation or yaw, zoom-in/out or scale, translation and out-of-plane rotations or pitch + roll. It is difficult to accurately

derive $\{\mathbf{R}_t^{t+1} \text{ and } \mathbf{T}_t^{t+1}\}$ from \mathbf{H}_t^{t+1} , and the errors in the solutions of the two components are coupled, i.e., an error in the translation estimate would induce a complementary error in the rotation estimate; this is highly undesirable.

Complementary sensors help mitigate this problem, and IMU is such a sensor which is present on quadrotors and can provide accurate angle measurements within a small interval. Thus, the problem of estimating cheap ego-motion reduces to finding \mathbf{T}_t^{t+1} (2D translation and zoom-in/out). Further, one can also obtain zoom-in/out using the altimeter on-board, but such an approach is noisy in a small interval. Hence we will estimate the 2D translation and zoom-in/out transformation which we refer to as ‘‘Pseudo-similarity’’ since it is one degree of freedom less than the similarity transformation (2D translation, zoom-in/out and yaw). We also call our network which estimates pseudo-similarity ‘‘PRGFlow’’. Note that, one can easily use our work to also estimate yaw without any added effort by changing the warping function. Mathematically, the pseudo-similarity transformation is given in Eq. 1, where W and H depict the image width and height, respectively.

$$\mathbf{x}_{t+1} = \begin{bmatrix} \frac{W}{2} & 0 & \frac{W}{2} \\ 0 & \frac{H}{2} & \frac{H}{2} \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1+s & 0 & t_x \\ 0 & 1+s & t_y \\ 0 & 0 & 1 \end{bmatrix} \mathbf{x}_t \quad (1)$$

We utilize the Inverse Compositional Spatial Transformer Networks (IC-STN) [39] for stacking multiple warping blocks for better performance. We extend the work in [39] to support pseudo-similarity and different warp types. A detailed study of which warp type performs the best is given in Sec. V-A. We further divide the experiments into table-top experiments (Sec. III) and flight experiments (Sec. IV).

III. TABLE-TOP EXPERIMENTS AND EVALUATION

A. Data Setup, Training and Testing Details

We train and test all our networks on the MS-COCO dataset [40], using the `train2014` and `test2014` splits for training and testing. During training, we obtain a random crop of size 300×300 px. (denoted as \mathcal{I}_t or \mathcal{I}_1), which is then warped using pseudo-similarity (2D translation + scale) to synthetically generate \mathcal{I}_{t+1} or \mathcal{I}_2 obtained by using a random warp parameter in the range $\gamma_1 = \pm [0.25 \ 0.20 \ 0.20]$ (unless specified otherwise). Then the center 128×128 px. patch are extracted (to avoid boundary effects) which are denoted as \mathcal{P}_1 and \mathcal{P}_2 , respectively. A stack of \mathcal{P}_1 and \mathcal{P}_2 patches of size $128 \times 128 \times 2N_c$ (N_c is the number of channels in \mathcal{P}_1 and \mathcal{P}_2 , i.e., 3 if RGB 1 if grayscale) is fed into the network to obtain the predicted warp parameters: $\tilde{\mathbf{h}} = [s \ t_x \ t_y]^T$.

Our networks were trained in PythonTM 2.7 using TensorFlow¹ 1.14 on a Desktop computer (described in Sec. III-E) running Ubuntu 18.04. We used a mini-batch-size of 32 for all the networks with a learning rate of 10^{-4} without any decay. We trained all our networks using the ADAM optimizer

for 100 epochs with early termination if we detect over-fitting on our validation set.

We tested our trained networks using two different configurations: (i) In-domain and (ii) Out-of-domain. In the in-domain testing (not to be confused with in-training dataset), we warped images from the `test2014` split of the MS-COCO dataset using a random warp parameter in range $\gamma_1 = \pm [0.25 \ 0.20 \ 0.20]$ (same warp range as training). This was used for evaluation as described in Sec. III-D. To commend on the generalization of the approach, we also tested it on out-of-domain warps, i.e., twice the warp range it was originally trained on, denoted by $\gamma_2 = \pm [0.50 \ 0.40 \ 0.40]$. This test was constructed to highlight the generalizability vs. speciality of the network configuration. Next we discuss the loss network architectures.

B. Network Architectures

We use five network architectures inspired by previous works. We call the networks as follows: VanillaNet [41], ResNet [42], ShuffleNet [43], MobileNet [44] and SqueezeNet [45]. Each network is composed of various blocks. The output of each block is used as an incremental warp as proposed in [39]. We use for each of these blocks the same name as the network, for e.g., VanillaNet has VanillaNet blocks. We use the following shorthand to denote the architecture. VanillaNet_a denotes VanillaNet with a VanillaNet blocks. We modify the networks from their original papers in the following manner: We exclude max-pooling blocks and replace them with stride in the previous convolutional block. Instead of varying sub-sampling rates (rates at which strides change with respect to depth of the network), we keep it fixed to the same value at every layer. All the architectures used are illustrated in Fig. 2. The architectures shown in Fig. 2 are repeated for every warping block used, where each block predicts the incremental warp as proposed in [39]. We exclude the number of filters for the purpose of clarity, but they can be found in the code provided with the supplementary material which will be released upon the acceptance for publication.

We also test two different sizes of networks, i.e., large (model size ≤ 8.3 MB) and small (Model size ≤ 0.83 MB). Here the model sizes are computed for storing `float32` values for each neuron weight. Due to file format packing there is generally a small overhead and the actual file sizes are slightly larger.

C. Loss Functions

The trivial way to learn the warp parameters directly is to use a supervised loss, since the labels are ‘‘free’’ as the data is synthetically generated on the fly. The loss function used in the supervised case is given as

$$\mathcal{L}_s = \underset{\tilde{\mathbf{h}}}{\operatorname{argmin}} \mathbb{E} \left(\|\tilde{\mathbf{h}} - \hat{\mathbf{h}}\|_2 \right), \quad (2)$$

where $\tilde{\mathbf{h}}$, $\hat{\mathbf{h}}$, \mathbb{E} are the predicted parameters, ideal parameters and expectation/averaging operator respectively. We also study the question ‘‘Does using self-supervised loss

¹<https://www.tensorflow.org/>



Figure 3. PRG Husky-360 γ platform used in flight experiments. (a) Top view, (b) front view, (c) down-facing leopard imaging camera.

give better out-of-domain generalization performance?” The unsupervised losses are generally complicated since they are under-constrained compared to the one in supervised approaches. Hence, we study different unsupervised loss functions. We can write the loss functions as:

$$\mathcal{L}_{u.s} = \underset{\tilde{\mathbf{h}}}{\operatorname{argmin}} \mathbb{E} \left(\mathcal{D} \left(\mathcal{W} \left(\mathcal{P}_1, \tilde{\mathbf{h}} \right), \mathcal{P}_2 \right) + \lambda_i \mathcal{R}_i \right), \quad (3)$$

where \mathcal{W} is a generic differentiable warp function, which can take on different mathematical formulations based on its second argument (model parameters), and \mathcal{D} represents a distance measuring image similarity between the image frames. Finally, λ_i , \mathcal{R}_i represent the i^{th} Lagrange multiplier and its corresponding regularization function. We experiment with different \mathcal{D} and \mathcal{R} functions described below.

$$\mathcal{D}_{L1} (A, B) = \mathbb{E} (\|A - B\|_1) \quad (4)$$

$$\mathcal{D}_{\text{Chab}} (A, B) = \left(\left((A - B)^2 + \epsilon^2 \right)^\alpha \right) \quad (5)$$

$$\mathcal{D}_{\text{SSIM}} (A, B) = \mathbb{E} \left(\frac{1 - \text{SSIM} (A, B)}{2} + \alpha (\|A - B\|_1) \right) \quad (6)$$

$$\mathcal{D}_{\text{Robust}} (A, B) = \mathbb{E} \left(\frac{b}{d} \left(\left(\frac{((A - B)/c)^2}{b} + 1 \right)^{d/2} - 1 \right) \right) \quad (7)$$

$$b = \|2 - \hat{\alpha}\|_1 + \epsilon; \quad d = \begin{cases} \hat{\alpha} + \epsilon & \text{if } \hat{\alpha} \geq 0 \\ \hat{\alpha} - \epsilon & \text{if } \hat{\alpha} < 0 \end{cases}$$

$$\hat{\alpha}_i = (2 - 2\epsilon_\alpha) \frac{e^{\alpha_i}}{e^{\alpha_i} + 1} \quad \forall i$$

Here, \mathcal{D}_{L1} is the generic l_1 photometric loss [46] commonly used for traditional images, $\mathcal{D}_{\text{Chab}}$ is the Chabonnier loss [47] commonly used for optical flow estimation, $\mathcal{D}_{\text{SSIM}}$ is the loss based on Structure Similarity [48] commonly used for learning ego-motion and depth from traditional image sequences and $\mathcal{D}_{\text{Robust}}$ is the robust loss function presented in [49]. Also, note that α has different connotation in each loss function.

The functions given above can take any generic input such as the raw image or a function of the image. In our paper we experiment with different inputs such as the raw RGB image, grayscale image, high-pass filtered image and

image cornerness score (denoted by \mathcal{I} , \mathcal{G} , $\mathcal{Z}(\mathcal{I})$ and $\mathcal{C}(\mathcal{I})$ respectively). The same set of functions can be used both as the metric function and the regularization function. We will denote the combination using a shorthand representation. Consider using the loss function with SSIM on raw images as the metric function and photometric L1 on high-pass filtered images as the regularization function with a Lagrange multiplier of 5.0, the shorthand for this function is given in Eq. 8.

$$\mathcal{D}_{\text{SSIM}} (\mathcal{I}) + 5.0 \mathcal{D}_{L1} (\mathcal{Z}(\mathcal{I})) \quad (8)$$

D. Evaluation Metrics

We use the following evaluation metrics to quantify the performance of each network. Let the predicted warp parameters be $\tilde{\mathbf{h}} = [\tilde{s} \ \tilde{t}_x \ \tilde{t}_y]^T$ and the ideal warp parameters be $\hat{\mathbf{h}} = [\hat{s} \ \hat{t}_x \ \hat{t}_y]^T$. W and H denote the image width and height respectively. Then the scale and translation error in pixel are given as:

$$\mathcal{E}_{\text{scale}} = \mathbb{M} \left(\sqrt{\frac{W^2 + H^2}{2}} |\tilde{s} - \hat{s}| \right) \quad (9)$$

$$\mathcal{E}_{\text{trans}} = \mathbb{M} \left(\sqrt{\frac{(W(\tilde{t}_x - \hat{t}_x))^2 + (H(\tilde{t}_y - \hat{t}_y))^2}{2}} \right) \quad (10)$$

Here, \mathbb{M} denotes the median value (we choose the median value over the mean to reject outlier samples with low texture).

We also convert errors to accuracy percentage as follows:

$$\mathcal{A} = \left(1 - \frac{\mathcal{E}_{\text{scale}} + \mathcal{E}_{\text{trans}}}{\mathbb{I}_{\text{scale}} + \mathbb{I}_{\text{scale}}} \right) \times 100\% \quad (11)$$

Here, $\mathbb{I}_{\text{scale}}$ and $\mathbb{I}_{\text{trans}}$ denote the identity errors for scale and translation respectively (error when the prediction values are zero).

²<http://nanopi.io/nanopi-neo-core2.html>

³<http://www.banana-pi.org/m2z.html>

⁴<https://coral.ai/products/dev-board/>

⁵<https://coral.ai/products/accelerator>

⁶<https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-tx2/>

⁷<https://up-board.org/>

⁸<https://www.asus.com/us/ROG-Republic-Of-Gamers/ROG-GL502VS/>

⁹<https://www.ibuypower.com/>

Table I
DIFFERENT COMPUTERS USED ON AERIAL ROBOTS.

Name	Cost (USD) ↓	Size ($l \times w \times h$ mm) ↓	Weight (g) ↓	CPU Arch.	CPU Clock Speed × Threads (GHz) ↑	RAM (GB) ↑	GPU (GOPs) ↑	Max. Power* (W) ↓	Ease of use ↑
NanoPi Neo Core 2 LTS ²	28	40 × 40 × 3	7	ARM	1.36 × 4	1.0	-	10	★★★★
BananaPi M2-Zero ³	24	65 × 30 × 5	15	ARM	1.00 × 4	1.0	-	10	★★★★
Google Coral Dev Board ⁴	150	88 × 60 × 22	136	ARM	1.50 × 4	1.0	4000	10	★★★★
Google Coral Accelerator ⁵	75	65 × 30 × 8	20	-	-	-	4000	4.5	★★★★
NVIDIA [®] Jetson [™] TX2 ⁶	600	87 × 50 × 48	200	ARM	2.00 × 6	8.0	1500	15	★★★★
Intel [®] Up Board ⁷	159	86 × 56 × 20	98	x86	1.92 × 4	1.0	-	10	★★★★★
Laptop ⁸	1600	391 × 267 × 31	2200	x86	3.40 × 8	16.0	6463	180	★★★★★

*Power consumption is for board not for chip.

E. Hardware Platforms

We tested multiple small factor computers and Single Board Computers (SBCs) (we’ll refer to both as computers or computing devices or hardware) which are commonly used on aerial platforms. The platforms tested are summarized in Table I. Details not present in the table are explained next.

We overclocked the NanoPi Neo Core 2 LTS to 1.368 GHz from the base clock of 1.08 GHz to obtain better performance. This necessitated the use of an active cooling solution for long duration operation (more than 3 mins continuously). The weight of the cooling solution (4 g) and micro-SD card (0.5 g) are not included in Table I. Without the active cooling solution the computer goes into a thermal shutdown which will be harmful during in-flight operation. One can change the clock speed to certain available frequencies between 0.48 GHz and 1.368 GHz. To run larger models on the NanoPi’s measly 1 GB of RAM we allocated 1 GB of SWAP memory on the Sandisk UHC-II class-10 microSD card which has much lower transfer speed as compared to RAM. The same cooling and SWAP solution were used for BananaPi M2-Zero as well. Note that, the Google Coral USB Accelerator is not officially supported with the NanoPi but we discovered a workaround which will be released in the accompanying supplementary material. Also, the Google Coral USB Accelerator is attached to a USB 2.0 port which limits its maximum performance when used with the NanoPi. The only reason why one would use BananaPi over the NanoPi is for the smaller width (30 mm versus 40 mm) which could be suited for a non X shaped quadrotor. However, the NanoPi is lighter, faster and has less area as compared to the BananaPi. Both NanoPi and BananaPi ran Ubuntu 16.04 LTS core with TensorFlow 1.14.

A significant speedup of upto 576× were obtained on the Coral Dev and the Coral USB accelerator when the original TensorFlow model was converted into TensorFlow-Lite and optimized for Edge TPU compilation. Without the Edge-TPU optimization the models run on the CPUs of these computers are far slower than the Tensor cores. To use both Coral Dev and Coral USB accelerator TensorFlow-Lite-Runtime 2.1 is used. The Coral Dev board ran Mendel Linux 1.5.

The NVIDIA[®] Jetson[™] TX2 in our setup is used with the Connect Tech’s Orbitty carrier board (weighing 41 g and its weight is included in Table I). This carrier board allows for the most compact setup with the TX2. Note that the NVIDIA[®] Jetson[™] TX2 can be used without the active heatsink for a short duration (less than 5 mins) reducing its weight by 59 g (a massive 30% reduction in weight without any loss of performance). However, extended operation without the

active heatsink can result in thermal shutdown or permanent damage to the computer. During our experiments, we set the operation mode to fix the CPU and GPU frequencies to the maximum available value, and this in-turn maxes out the power consumption (the steps to achieve this will be released in the accompanying supplementary material). The TX2 ran Linux For Tegra (L4T) R28.2.1 installed using Jetpack 3.4 with TensorFlow 1.11.

We use neither the AVX (not supported on hardware) nor the SSE (not supported by the TensorFlow version supported on Up board) instruction set on the Intel[®] Up board. We speculate huge speed-ups if a future version of TensorFlow supports SSE on the Up board. The Up board ran Ubuntu 16.04 LTS with TensorFlow 1.11.

The laptop is an Asus ROG GL502VS with Intel[®] Core[™] i7-6700HQ and GTX 1070 GPU and weighs about 2200 g which can be used on a large (≥ 650 mm sized) quadrotor. Also, note that an Intel NUC coupled to an eGPU can also be used with similar specifications but this setup will still be heavier, arguably more expensive and probably less reliable than gutting out a gaming laptop. The laptop ran Ubuntu 18.04 LTS with TensorFlow 1.14.

The desktop PC is a custom built full tower PC from iBUYPOWER with an Intel[®] Core[™] i9-9900KF and NVIDIA[®] Titan-Xp GPU which cannot be used on an aerial robot but is included to serve as a reference. The desktop PC ran Ubuntu 18.04 LTS with TensorFlow 1.14. Whenever possible we use the NEON SIMD instruction set for ARM computers and the SSE instructions for x86 systems.

IV. FLIGHT EXPERIMENTS AND EVALUATION

This section presents real-world experiments flying various simple trajectories with odometry estimation using PRGFlow. We use the lowest avg. pixel error 8.3 MB (large) model for evaluation (ResNet₄ with T×2, S×2 warping configuration). The network outputs a $\mathbb{R}^{3 \times 1}$ vector denoted as $\tilde{\mathbf{h}} = [\tilde{s} \quad \tilde{t}_x \quad \tilde{t}_y]^T$. The predicted translational pixel velocities (global optical flow) $[\tilde{t}_x \quad \tilde{t}_y]^T$ are converted to real world velocities by scaling them using the focal length and the depth (adjusted altitude) [14]. A similar treatment is given to the Z-pixel velocity s . Also, we obtain attitude using a Madgwick filter [50] from a 9-DoF IMU, which is used to remove rotation between consecutive frames, and feed the values into the network. Finally, to obtain odometry, we simply perform dead-reckoning on the velocities obtained by our network. *Our networks were trained on MS-COCO as described in III-A, and*

Table II
QUANTITATIVE EVALUATION OF DIFFERENT WARPING COMBINATION FOR PSEUDO-SIMILARITY ESTIMATION.

Network (Warping)	$\mathcal{E}_{\text{scale}}$ (px.)		$\mathcal{E}_{\text{trans}}$ (px.)		FLOPs (G) ↓	Num. Params (M) ↓
	γ_1	γ_2	↓	↓		
Identity	11.4	22.8	10.3	20.4	–	–
VanillaNet ₁ (PS×1)	2.4	15.0	1.3	12.5	0.37	2.07
VanillaNet ₁ (PS×1) DA	4.1	17.7	2.3	14.2	0.37	2.07
VanillaNet ₂ (PS×2)	2.2	9.9	1.4	12.4	0.42	2.17
VanillaNet ₂ (S×1, T×1)	2.5	15.2	1.5	12.2	0.46	2.10
VanillaNet ₂ (T×1, S×1)	2.5	15.1	1.5	12.5	0.42	2.15
VanillaNet ₄ (PS×4)	2.3	11.9	1.5	14.9	0.42	2.15
VanillaNet ₄ (S×2, T×2)	2.6	15.4	1.6	12.6	0.46	2.08
VanillaNet ₄ (T×2, S×2)	2.0	8.5	1.5	12.5	0.46	2.08
VanillaNet ₄ (T×2, S×2) γ_2	2.7	2.8	4.6	7.2	0.46	2.08

Table III
QUANTITATIVE EVALUATION OF DIFFERENT NETWORK ARCHITECTURES FOR PSEUDO-SIMILARITY ESTIMATION USING T×2, S×2 WARPING BLOCK FOR LARGE MODEL (≤8.3 MB).

Network	$\mathcal{E}_{\text{scale}}$ (px.)		$\mathcal{E}_{\text{trans}}$ (px.)		FLOPs (G) ↓	Num. Params (M) ↓
	γ_1	γ_2	↓	↓		
Identity	11.4	22.8	10.3	20.4	–	–
VanillaNet ₄	1.9	6.4	1.5	12.4	0.46	2.08
ResNet ₄	1.7	15.1	0.9	10.1	0.59	2.12
SqueezeNet ₄	2.1	5.7	2.2	13.8	2.20	2.12
MobileNet ₄	4.0	14.2	1.6	12.0	0.41	2.04
ShuffleNet ₄	6.4	17.4	3.0	13.9	1.20	2.10

are used in the flight experiments without any fine-tuning or re-training to highlight the generalizability of our approach. Note that the performance can be significantly boosted by multi-frame fusion of predictions, and we leave this for future work.

A. Experimental Setup

We tested our algorithm on the PRG Husky-360 γ platform⁹: a modified version of the Parrot[®] Bebop 2 for its ease of use which was originally designed for pedagogical reasons. It is equipped with a down facing Leopard Imaging LI-USB30-M021 global shutter camera¹⁰ with a 16 mm lens which gives a diagonal field of view of $\sim 22^\circ$ (Refer to Fig. 3). The PRG Husky also contains an 9-DoF IMU and a down-facing sonar for attitude and altitude measurements respectively. Higher level control commands are given by an on-board companion computer: Intel[®] Up Board at 20 Hz. The images are recorded at 90 Hz at a resolution of 640 × 480 px. The overall takeoff weight of the flight setup is 730 g (which gives a thrust-to-weight ratio of ~ 1.5) with diagonal motor to motor dimension of 360 mm.

The flight experiments were conducted in the Autonomy Robotics and Cognition (ARC) lab’s netted indoor flying space at the University of Maryland, College Park. The total flying volume is about $6 \times 5.5 \times 3.5 \text{ m}^3$. A Vicon motion capture system with 12 vantage V8 cameras were used to obtain ground truth at 100 Hz.

The PRG Husky was tested on five trajectories: circle, moon, line, figure8 and square which involve change in both attitude and altitude with an average velocity of about 0.5 ms^{-1} and a maximum velocity of 1.5 ms^{-1} .

⁹<https://github.com/prgumd/PRGFlyt/wiki/PRGHusky>

¹⁰<https://leopardimaging.com/product/usb30-cameras/usb30-camera-modules/li-usb30-m021/>

Table IV
QUANTITATIVE EVALUATION OF DIFFERENT NETWORK ARCHITECTURES FOR PSEUDO-SIMILARITY ESTIMATION USING T×2, S×2 WARPING BLOCK FOR SMALL MODEL (≤0.83 MB).

Network	$\mathcal{E}_{\text{scale}}$ (px.)		$\mathcal{E}_{\text{trans}}$ (px.)		FLOPs (G) ↓	Num. Params (M) ↓
	γ_1	γ_2	↓	↓		
Identity	11.4	22.8	10.3	20.4	–	–
VanillaNet ₄	3.3	8.9	3.1	14.0	0.18	0.21
ResNet ₄	4.4	12.5	2.4	12.1	0.20	0.20
SqueezeNet ₄	2.4	5.6	4.0	14.9	0.19	0.20
MobileNet ₄	8.3	18.7	3.7	13.4	0.16	0.20
ShuffleNet ₄	8.3	17.6	4.6	15.7	0.13	0.21

Table V
QUANTITATIVE EVALUATION OF DIFFERENT NETWORK INPUTS FOR PSEUDO-SIMILARITY ESTIMATION USING T×2, S×2 WARPING BLOCK FOR LARGE MODEL (≤8.3 MB).

Testing Data (Training Data)	$\mathcal{E}_{\text{scale}}$ (px.)		$\mathcal{E}_{\text{trans}}$ (px.)	
	γ_1	γ_2	↓	↓
Identity	11.4	22.8	10.3	20.4
\mathcal{I} (\mathcal{I})	2.0	8.5	1.5	12.5
\mathcal{G} (\mathcal{I})	1.8	6.3	1.5	12.3
\mathcal{G} (\mathcal{G})	2.7	14.1	1.6	12.7
\mathcal{Z} (\mathcal{I}) (\mathcal{Z} (\mathcal{I}))	13.1	9.4	10.4	16.0
\mathcal{G} (\mathcal{Z} (\mathcal{I}))	11.8	20.7	9.8	19.8
\mathcal{I} (\mathcal{Z} (\mathcal{I}))	13.1	22.5	10.5	20.1
\mathcal{Z} (\mathcal{I}) (\mathcal{G})	8.5	19.8	4.1	17.6
\mathcal{Z} (\mathcal{I}) (\mathcal{Z})	17.2	20.1	4.2	17.4

Table VI
QUANTITATIVE EVALUATION OF DIFFERENT LOSS FUNCTIONS FOR PSEUDO-SIMILARITY ESTIMATION USING PS×1 WARPING BLOCK FOR LARGE MODEL (≤8.3 MB).

Loss Function (Architecture)	$\mathcal{E}_{\text{scale}}$ (px.)		$\mathcal{E}_{\text{trans}}$ (px.)		FLOPs (G) ↓	Num. Params (M) ↓
	γ_1	γ_2	↓	↓		
Identity	11.4	22.8	10.3	20.4	–	–
Supervised \mathcal{L}_s (VanillaNet ₁)	2.4	15.0	1.3	12.5	0.37	2.07
$\mathcal{D}_{\text{Robust}}$ (\mathcal{I} , \mathcal{C} (\mathcal{I})) (ResSqueezeNet ₁)	12.9	25.2	7.2	11.7	1.01	2.18
$\mathcal{D}_{\text{SSIM}}$ (\mathcal{I}) (ResSqueezeNet ₁)	3.4	21.2	6.0	13.8	1.01	2.18
$\mathcal{D}_{\text{SSIM}}$ (\mathcal{I}) + 0.1 $\mathcal{D}_{\mathcal{L}_1}$ (\mathcal{C} (\mathcal{I})) (ResSqueezeNet ₁)	2.0	16.1	6.2	14.6	1.01	2.18
$\mathcal{D}_{\text{SSIM}}$ (\mathcal{I}) + 0.1 $\mathcal{D}_{\mathcal{L}_1}$ (\mathcal{Z} (\mathcal{I})) (ResSqueezeNet ₁)	2.7	16.6	6.4	13.6	1.01	2.18
$\mathcal{D}_{\mathcal{L}_1}$ (DB (\mathcal{E})) [13]	5.4	17.7	3.7	16.5	4.92	3.6
$\mathcal{D}_{\text{Chab}}$ (DB (\mathcal{E})) [13]	5.1	17.1	3.4	16.7	4.92	3.6
Supervised DB (\mathcal{E}) [13]	4.1	16.2	3.3	15.1	4.92	3.6

Table VII
QUANTITATIVE EVALUATION OF DIFFERENT COMPRESSION METHODS FOR PSEUDO-SIMILARITY ESTIMATION USING PS×1 WARPING.

Method	$\mathcal{E}_{\text{scale}}$ (px.)		$\mathcal{E}_{\text{trans}}$ (px.)		$\mathcal{E}_{\text{scale}}$ (DA) (px.)		$\mathcal{E}_{\text{trans}}$ (DA) (px.)	
	γ_1	γ_2	↓	↓	γ_1	γ_2	↓	↓
Identity	11.4	22.8	10.3	20.4	11.4	22.8	10.3	20.4
Teacher from Scratch	2.4	15.0	1.3	12.5	4.3	17.2	2.5	14.1
Student from Scratch	3.5	9.9	2.8	13.2	4.2	16.8	4.2	16.0
Projection Loss Student [51]	3.7	10.9	2.8	13.1	8.0	17.7	4.2	15.2
Model Distillation Student [52]	3.8	12.3	2.7	13.4	7.2	17.7	4.0	15.2

Table VIII
COMPARISON OF PRGFLOW WITH DIFFERENT CLASSICAL METHODS.

Method	$\mathcal{E}_{\text{scale}}$ (px.) ↓	$\mathcal{E}_{\text{trans}}$ (px.) ↓	$\mathcal{E}_{\text{scale}}$ (DA) (px.) ↓	$\mathcal{E}_{\text{trans}}$ (DA) (px.) ↓	Time (ms) ↓
Identity	11.4	10.3	11.4	10.3	–
Supervised \mathcal{L}_s	1.9	1.5	4.1	2.3	1.1*
FFT Alignment [53]	0.3	0.1	13.4	6.0	35.5
SURF [54]	0.4	0.1	11.2	1.3	17.6
ORB [55]	0.6	0.1	11.5	1.4	12.9
FAST [56]	0.8	0.2	12.0	1.3	55.9
Brisk [57]	0.7	0.2	13.0	1.4	38.7
Harris [58]	0.4	0.1	10.9	1.4	60.2

* On all cores.

V. DISCUSSION

A. Algorithmic Design

We answer the following questions in this section:

- 1) What is the best warp combination?
- 2) What network architecture is the best? Does the best network architecture vary with respect to the number of

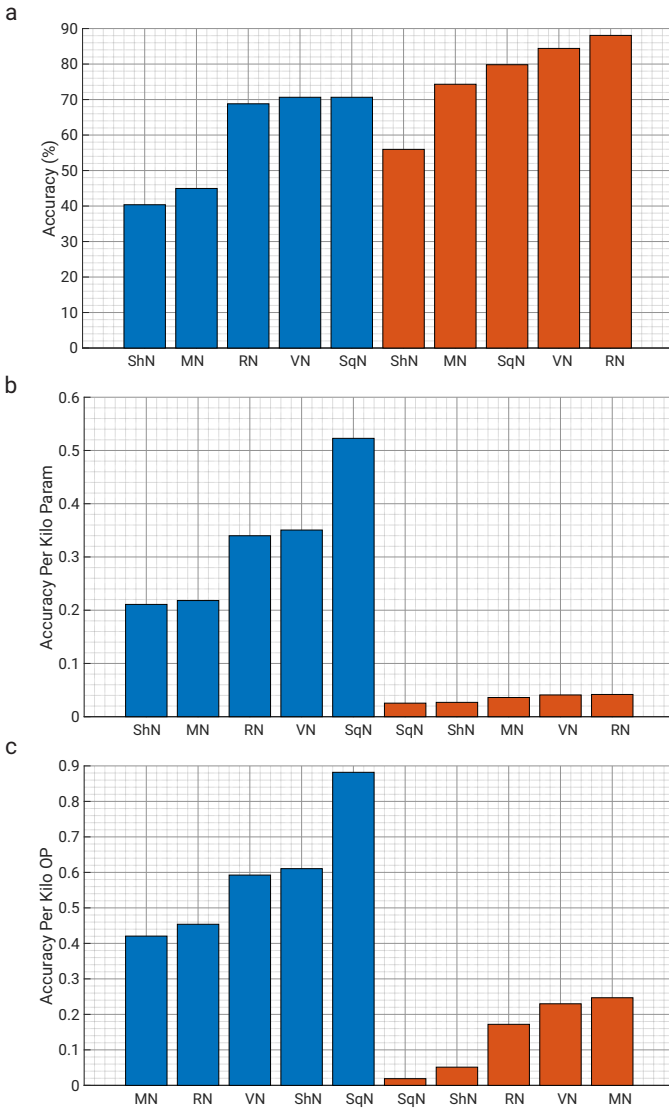


Figure 4. (a) Accuracy, (b) Accuracy per Kilo param, (c) Accuracy per Kilo OP for different network architectures. Blue and orange histograms denote small (≤ 0.83 MB) and large (≤ 8.3 MB) networks respectively. Here the following shorthand is used for network names: VN: VanillaNet, RN: ResNet, SqN: SqueezeNet, MN: MobileNet and ShN: ShuffleNet. All networks use $T \times 2, S \times 2$ warping configuration.

parameters?

- 3) How does the input affect performance of a network?
- 4) What is the best way to compress a model?
- 5) When is it advisable to choose deep learning for odometry estimation?

The performance of different IC-STN warp combinations for VanillaNet are given in Table II. The in-domain and out-of-domain test results have a headings of γ_1 and γ_2 respectively. All the networks were trained using supervised l_2 loss (Eq. 2). Note that, the networks were constrained to be ≤ 8.3 MB. Under this condition, one can clearly observe the following trend when only using PS blocks for warping – as the number of warping blocks increases the performance reaches a maximum and then deteriorates. This is because the number of neurons per warp block directly affects the performance, and increasing the number of warp blocks

without increasing model size hurts accuracy when the number of neurons per warp block become small. An interesting observation is that predicting translation before scale almost always results in better performance and also decoupling the predictions of scale and translation using S and T blocks generally results in better performance (lower pixel error) as compared to predicting them together using PS blocks. This is serendipitous. As drift in the X and Y position are generally higher, one could obtain these positions faster (since only a part of the network is used for output) and at higher frequency than the Z position [59].

We also observe that training and testing with data augmentation of brightness, contrast, hue, saturation and Gaussian noise worsens the average performance by 73% (comparing rows 2 and 3 in Table II) since the network now has to learn to be agnostic to a myriad of image changes. Also, training on γ_2 (last two rows of Table II) decreases the average error by 62% when testing on γ_2 - we speculate it will further decrease with increasing model size. This shows that as the deviation range in training increases, one can require more parameters for the same average performance. *Overall the best performing warp configuration is $T \times 2, S \times 2$ when considering an average of both in-domain and out-of-domain tests.*

Next, let us study the performance with different network architectures. We use $T \times 2, S \times 2$ warps for all the networks, and they are trained using supervised l_2 loss (Eq. 2). The results for networks constrained by model size ≤ 8.3 MB and ≤ 0.83 MB can be found in Tables III and IV respectively. For ease of analysis, the results are also visually depicted in Fig. 4. One can clearly observe that ResNet gives the best performance for both small and large networks (Fig. 4a) and should be the network architecture of choice when designing for maximizing accuracy without any regard to the number of parameters or amount of OPs (operations). However, if one has to prioritize maximizing accuracy whilst minimizing number of parameters, then SqueezeNet and ResNet would be the choice for smaller and larger networks respectively (Fig. 4b). Another trend one can observe is that one would need $10 \times$ more parameters for a 19% increase in accuracy. Lastly, if one has to prioritize in maximizing accuracy whilst minimizing the number of OPs, then SqueezeNet and MobileNet would be the choice for smaller and larger networks respectively (Fig. 4c). Clearly, the most optimal architecture in-terms of accuracy, number of parameters and OPs is SqueezeNet for smaller networks and ResNet for larger networks. *The effect of the choice of network architecture is fairly significant on the accuracy, number of parameters and OPs and needs to be carefully considered when designing a network for deployment on an aerial robot.* Also, note that sometimes using a classical approach to solve a small part of the problem can significantly simplify the learning problem for the network thereby maximizing accuracy and minimizing the number of parameters and OPs [1].

To gather more insight into what data representation is more important for odometry estimation, we explore training and testing on the following data representations: (a) RGB image, (b) Grayscale image, (c) High pass filtered

image. We choose the $T \times 2$, $S \times 2$ warp configuration and the VanillaNet₄ architecture trained using supervised l_2 loss constrained by model size ≤ 8.3 MB. Table V summarizes the results obtained. Surprisingly, training on RGB images and evaluating on RGB images gives worse performance in the testing both for in-domain (γ_1) and out-of-domain (γ_2) ranges than testing on grayscale data. We speculate that this is due to conflicting information in multiple channels. Another surprising observation is that training and/or testing on high-pass filtered images results in large errors, which is contrary to the classical approaches. We speculate that this is because conventional neural networks rely on “staticity” of the image pixels (image pixels change slowly and are generally smooth).

We also explore the state-of-the-art self-supervised/unsupervised loss functions to test for claims of better out-of-domain generalization. We choose half-sized (≤ 4.15 MB) ResNet (to output scale) and SqueezeNet (to output translation) denoted as ResSqueezeNet trained using $PS \times 1$ blocks for this experiment since it empirically gave us the best results (we exclude other architecture results for the purpose of brevity). We also include results from a VanillaNet₁ trained using $PS \times 1$ blocks using the supervised l_2 loss function to serve as a reference (Refer to Table VI). We observe that scale error when using SSIM and cornerness (obtained using a heatmap as the output of [60]) in the loss approaches the performance of the supervised network, but the translation error is almost three times that of the supervised network. Surprisingly, the supervised network also performs better than most unsupervised networks on out-of-domain tests. *This hints that we need better loss functions for unsupervised methods and better network architectures to take advantage of these unsupervised losses.* From a practitioner’s point of view, the supervised networks perform better and generalize better to out-of-domain. Another keen insight is that focusing on crafting better supervised loss functions may lead to a massive boost in performance.

Finally, we explore different strategies to compress the network. We specifically consider a setup where we compress a 8.3 MB model VanillaNet₁ $PS \times 1$ to a 0.83 MB model VanillaNet₁ $PS \times 1$. We test three different methods, (a) direct dropping of weights, (b) projection inspired loss and (c) model distillation. In the first method, we reduce the number of neurons and number of blocks to reduce the number of weights and train the smaller network from scratch. In the second method, we train both the larger and smaller networks together as given by Eq. 12 [51]. Here, \hat{h}_T, \hat{h}_S denote the predictions from the teacher and student network respectively. We choose $\lambda_1, \lambda_2, \lambda_3$ as 1.0, 1.0 and 0.1 respectively. In the third method, we use an already trained teacher network (large 8.3 MB model) and define the loss to learn the predictions of the teacher using the student (small 0.83 MB model) as given by Eq. 13 [52]. When no data augmentation is used, we observe that directly dropping weights gives the best performance (Refer to Table VII). However, when the data augmentation is added the model distillation gives the best results, albeit only slightly better than directly dropping of weights. This observation is contrary to that observed with classification

networks where massive boosts in performance are observed when using either Eq. 12 or 13. From a practitioner’s point of view, the simplest method of directly dropping weights work the best for regression networks like the one used in this work and can provide up to $10 \times$ savings in the number of parameters and OPs at the cost of $\sim 19\%$ accuracy. The same effect is observed when training and testing with and without data augmentation.

$$\mathcal{L}_{\text{Proj}} = \lambda_1 \mathcal{L}_s(\hat{h}, \widetilde{h}_T) + \lambda_2 \mathcal{L}_s(\hat{h}, \widetilde{h}_S) + \lambda_3 \mathcal{L}_s(\widetilde{h}_T, \widetilde{h}_S) \quad (12)$$

$$\mathcal{L}_{\text{TS}} = \mathcal{L}_s(\widetilde{h}_T, \widetilde{h}_S) \quad (13)$$

To address the elephant in the room, we try to answer the following question: “When should one use deep learning over a traditional approach?” We compare the proposed deep learning approach PRGFlow to the common method of fast feature matching on aerial robots in Table VIII. Note that the runtime for traditional methods are given for MATLAB¹¹ implementations on one thread of the Desktop PC to standardize the libraries and optimizations used. Up to $5 \times$ and $10 \times$ speedup can be obtained using efficient C++ implementations on a multi-threaded CPU and GPU respectively (we don’t explore this in our work). One can clearly observe that even with C++ implementations the traditional handcrafted features are slower than the deep learning methods which can utilize the parallel hardware accelerations on GPUs. Though on the surface it seems like the traditional methods give far superior performance in terms of accuracy (lower error), the efficacy of deep learning approaches are brought into limelight when we train and test with data augmentations. This simulates a bad quality camera common on smaller aerial robots. The drop in accuracy (from no-noise data to noisy data) in deep learning approaches is much less than that compared to traditional approaches, i.e., *deep learning approaches on an average fail more gracefully compared to traditional approaches.*

B. Hardware Aware Design

We answer the following questions in this section:

- 1) How fast do different network architectures run on a variety of computing platforms subject to weight and volume constraints?
- 2) What is the most efficient network architecture for a given hardware abiding the SWAP constraints?
- 3) How does varying network width and depth affect the speed of different network architectures on different computing platforms?
- 4) Which hardware setup is more power efficient?
- 5) How significant is power used for computing compared to power used by other quadrotor components?

In the wise words of Alan Kay, a pioneer in computer engineering “*People who are really serious about software should make their own hardware*” one would ideally want

¹¹<https://www.mathworks.com/products/matlab.html>

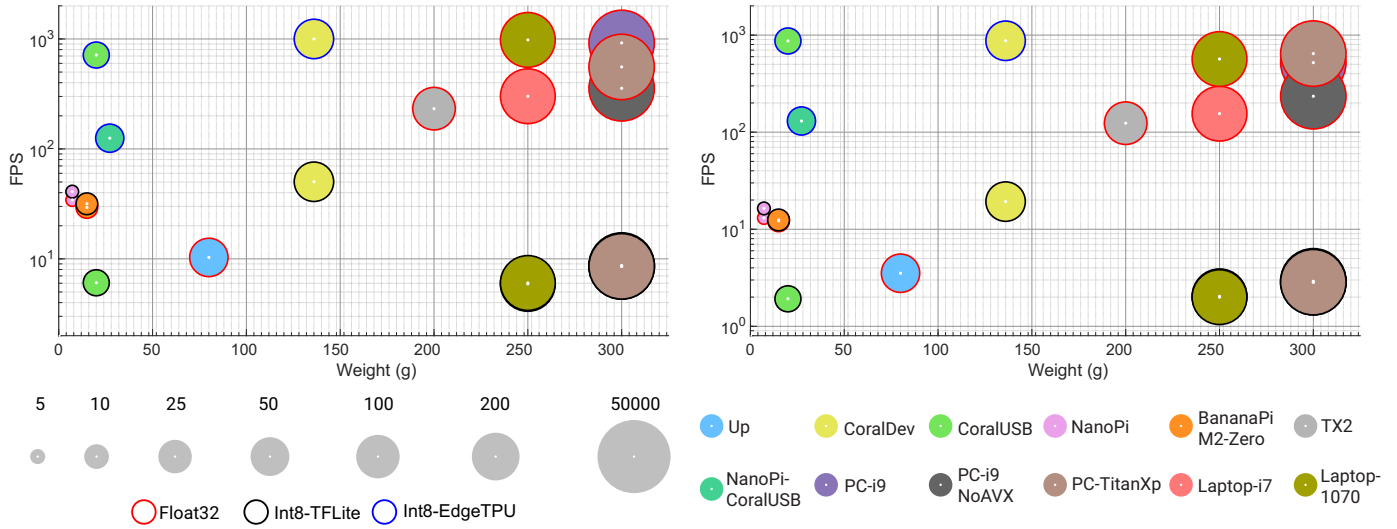


Figure 5. Weight vs. FPS for VanillaNet₄ (T×2, S×2) on different hardware and software optimization combinations. Left: small (≤ 0.83 MB) model, right: large (≤ 8.3 MB) model. The radius of each circle is proportional to log of volume of each hardware (this is shown in the legend below the plots with volume indicated on top of each legend in cm^3). The outline on each sample indicates the configuration of quantization or optimization used (Float32 (red outline) is the original TensorFlow model without any quantization or optimization, Int8-TFLite (black outline) is the TensorFlow-Lite model with 8-bit Integer quantization and Int8-EdgeTPU (blue outline) is the TensorFlow-Lite model with 8-bit Integer quantization and Edge-TPU optimization). The samples are color coded to indicate the computer it was run on (shown in the legend on the bottom). Also note that, Laptop and PC (Desktop) weight and volume values are not to actual scale for visual clarity in all images. All the figures in this paper use the same legend and color coding for ease of readability.

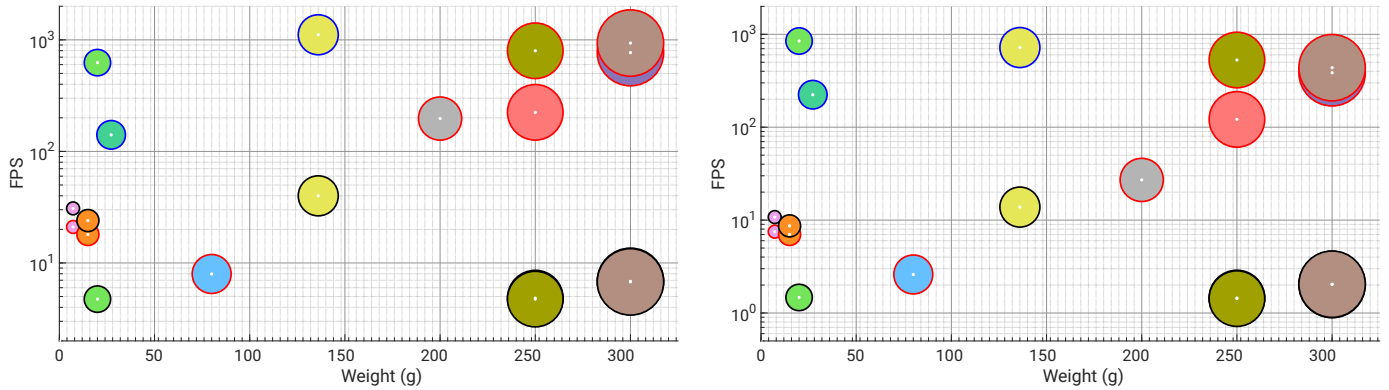


Figure 6. Weight vs. FPS for ResNet₄ (T×2, S×2) on different hardware and software optimization combinations. Left: small (≤ 0.83 MB) model, right: large (≤ 8.3 MB) model. The radius of each circle is proportional to log of volume of each hardware.

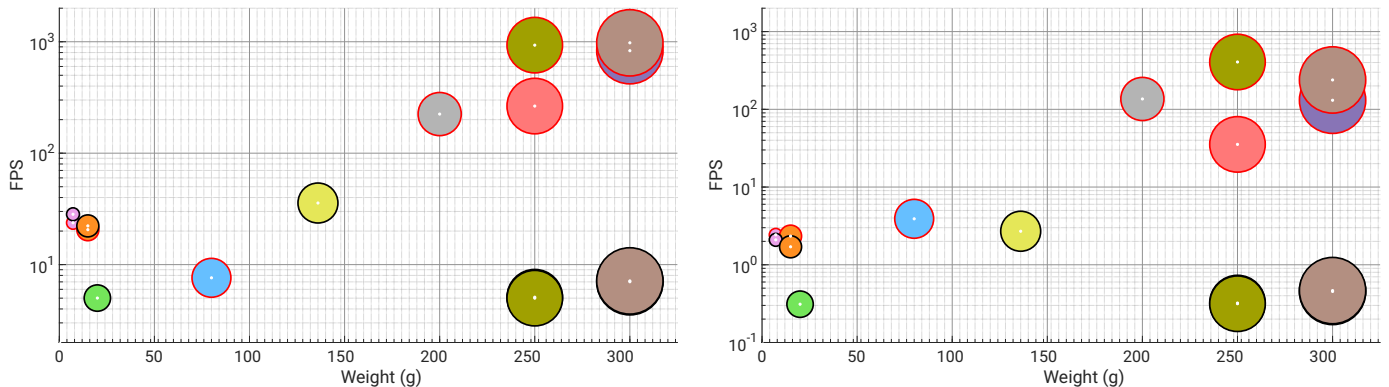


Figure 7. Weight vs. FPS for SqueezeNet₄ (T×2, S×2) on different hardware and software optimization combinations. Left: small (≤ 0.83 MB) model, right: large (≤ 8.3 MB) model. The radius of each circle is proportional to log of volume of each hardware.

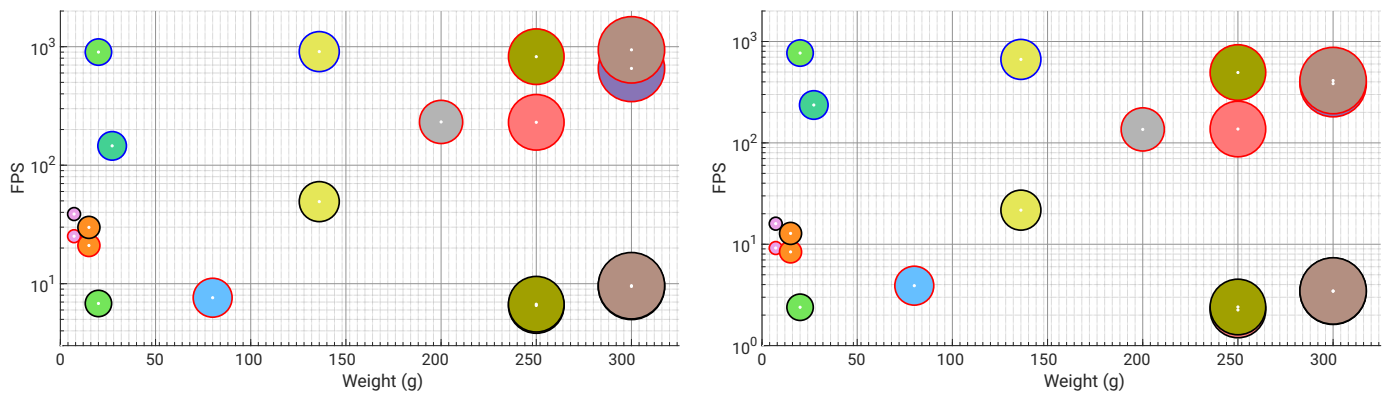


Figure 8. Weight vs. FPS for MobileNet₄ (T×2, S×2) on different hardware and software optimization combinations. Left: small (≤ 0.83 MB) model, right: large (≤ 8.3 MB) model. The radius of each circle is proportional to log of volume of each hardware.

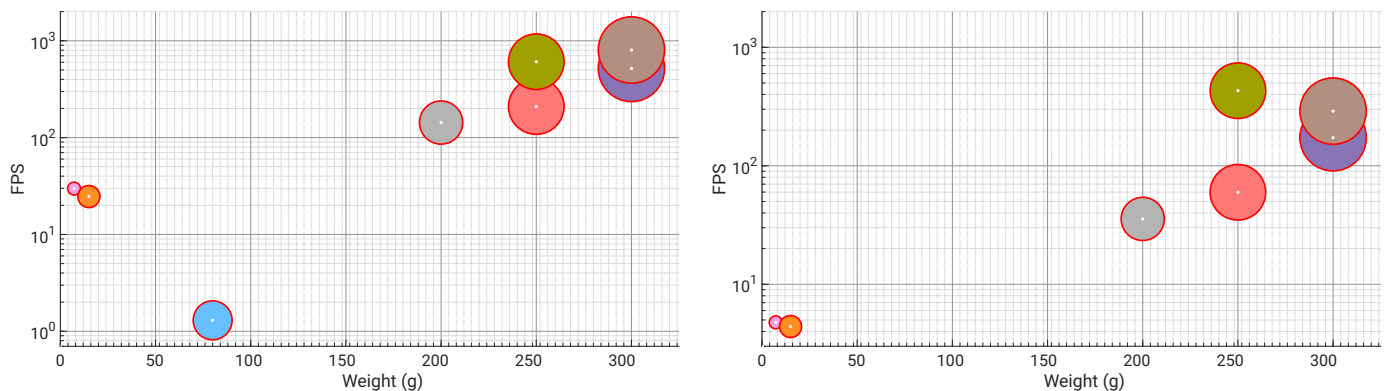


Figure 9. Weight vs. FPS for ShuffleNet₄ (T×2, S×2) on different hardware and software optimization combinations. Left: small (≤ 0.83 MB) model, right: large (≤ 8.3 MB) model. The radius of each circle is proportional to log of volume of each hardware.

Table IX
DIFFERENT-SIZED QUADROTOR CONFIGURATION WITH RESPECTIVE COMPUTERS.

Quadrotor Size (mm)	Propeller Size (mm)	Motor	Computing Board	Computer Weight (g) ↓	Total Weight (g) ↓	Auto. Thrust to Weight ↑	Auto. Hover Power (W) ↓	Hover Power (W) ↓	ρ ↓
75	40	Happymodel SE0706 15000KV	Nano Pi	7	62	1.43	33	17	1.92
120	63	T-Motor F15 Pro KV6000	Nano Pi + Coral USB	29	209	4.67	98	72	1.36
160	76	T-Motor F20II KV3750	Nano Pi + Coral USB	29	279	4.93	75	49	1.53
210	152	EMAX RS2306 KV2750	Coral Dev Board	136	536	9.91	106	64	1.67
360	178	T-Motor F80 Pro KV1900	NVIDIA® NVIDIA® Jetson™ TX2	200	1100	7.69	318	222	1.43
500	254	iFlight XING 2814 1100KV	NVIDIA® NVIDIA® Jetson™ Xavier AGX	600	2000	4.98	657	551	1.19
650	381	Tarot 4414 KV320	Intel® NUC + NVIDIA® Jetson™ Xavier AGX	1300	3900	2.72	603	405	1.49

Table X
TRAJECTORY EVALUATION FOR FLIGHT EXPERIMENTNTS OF PRGFLOW.

Trajectory	Error (m) ↓				Traj. Length (m)
	X	Y	Z	RMSE	
Circle	0.04	0.06	0.14	0.12	12.21
Moon	0.06	0.08	0.08	0.09	11.67
Line	0.06	0.06	0.10	0.09	3.84
Figure8	0.03	0.03	0.05	0.05	10.91
Square	0.04	0.05	0.10	0.08	10.77

to design a hardware chip for a specific SWAP constraint dictated by the size and amount of features desired in the aerial robot. This can generally only be achieved by the elite drone manufacturing companies owing to the exorbitant non-recurring engineering cost, thereby, putting research labs at a handicap. However, due to the rising Internet of Things (IoT) technologies and computers required to fit tight SWAP constraints researchers can repurpose these computers for

efficient utilization on quadrotors (or aerial robots in general). In this spirit, we limit our analysis to commonly used computers designed for IoT purposes which are repurposed for use on aerial robots. The computers used in our experiments are summarized in Table I and discussed in more detail in Subsec. III-E.

Refer to Figs. 5, 6, 7, 8 and 9 for a plot of Weight vs. FPS (Frames Per Second) vs. volume of the computer for VanillaNet, ResNet, SqueezeNet, MobileNet and SqueezeNet respectively. All the networks include both small (≤ 0.83 MB) and large (≤ 8.3 MB) configurations training using supervised l_2 loss function and optimized using different post-quantization optimizations such as Int8-TF Lite and Int8-EdgeTPU. We exclude results from Float32-TF Lite due to inferior performance without any significant speedups as compared to the original Float32 model. One can clearly observe that for all the networks the

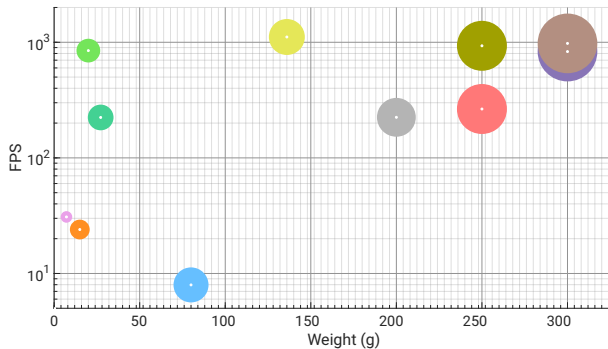


Figure 10. Weight vs. FPS for the best model architecture on each hardware coupled to the best software optimization combination. The radius of each circle is proportional to log of volume of each hardware. The best model architecture and model optimization for each hardware are: Up: ResNet_S-Float32, CoralDev: ResNet_S-Int8-EdgeTPU, CoralUSB: ResNet_S-Int8-EdgeTPU, NanoPi: ResNet_S-Int8, BananaPiM2-Zero: ResNet_S-Int8, TX2: SqueezeNet_S-Float32, Laptop-i7: SqueezeNet_S-Float32, Laptop-1070: SqueezeNet_S-Float32, PC-i9: SqueezeNet_S-Float32, PC-TitanXp: SqueezeNet_S-Float32. All networks use T×2, S×2 configuration and *S* and *L* subscripts indicate small and large networks respectively. The best network for each hardware was chosen with the avg. error ≤ 2.5 px. and the configuration which gives the highest FPS.

desktop and laptop give the best speed (highest FPS) but also have the largest weights. It would be advisable to gut out a gaming laptop and use it on a larger quadrotor (≥ 650 mm) due to the availability of integrated NVIDIA[®] mobile GPUs with a large amount of CUDA[®] cores which can be used to accelerate both deep learning and traditional computer vision tasks. An important factor to realize is that using Int8-EdgeTPU optimization to run on either the Coral Dev board or the Coral USB accelerator can provide significant speed-ups of up to 52× compared to Int8-TFLite without significant loss in accuracy. However, not all operations are supported in TensorFlow Lite and even less operations are supported for EdgeTPU optimization. Also, a drop in speed when going from a smaller model to a larger model is less significant in coral boards due to efficient TPU architecture. Hence, it is advisable to use the Coral Dev board or the Coral USB accelerator whenever possible. We also exclude Intel[®]'s Movidius Neural Compute sticks in our analysis since they provide inferior performance than Coral boards, are harder to use, weigh more and are larger.

A non-obvious observation is that the Int8-TFLite execution speed is much lower than the Float32 model on laptops and desktops (both CPU and GPU). This is because of lack of optimized 8-bit integer instruction sets which are generally present in lower-end ARM computers such as the NanoPi and BananaPi. We can also observe a similar drop in performance of up to 2.6× when AVX and SSE optimized instruction sets are not used on the desktop (Fig. 5, datapoint indicated as PC-i9 NoAVX). The lack of good performance of the Up board can also be pinned to the lack of AVX and SSE instruction sets and should be avoided for neural network tasks if not coupled to a neural network accelerator such as the Intel[®] Movidius compute stick or the Coral USB accelerator. We also

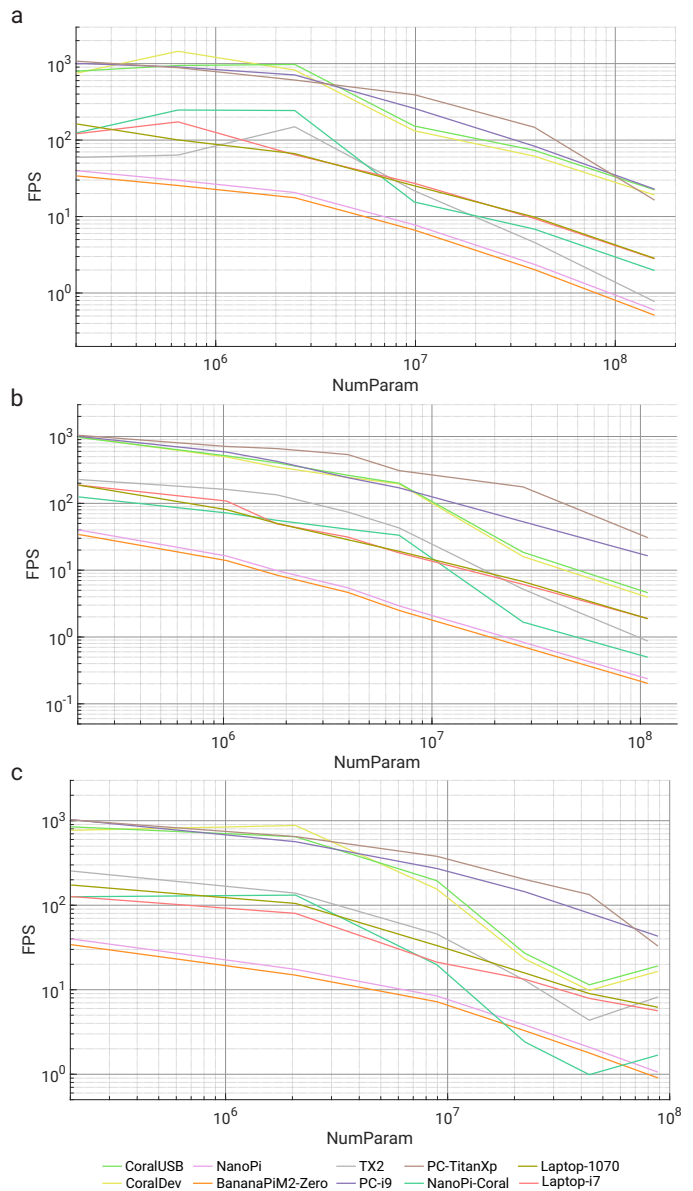


Figure 11. Num. of Params vs. FPS (a) when only increasing the depth of the network while keeping width constant, (b) when only increasing the width of the network while keeping depth constant, (c) when increasing a combination of depth and width of the network for different computers.

observe speedups of upto 1.7× on the NanoPi and BananaPi by converting a Float32 model to Int8-TFLite due to accelerated NEON instruction sets. Also, note that ShuffleNet models do not support Int8-TFLite and Int8-EdgeTPU optimizations due to unsupported layers. Another interesting observation is that, ResNet (both smaller and larger) and smaller SqueezeNet models achieve almost the same speed on both CPUs and GPUs on the desktop computer.

We choose as the best network architecture and configuration (small versus large) for each hardware, the one which has an average error ≤ 2.5 px. and gives the maximum speed. These results are illustrated in Fig. 10. For smaller boards like NanoPi and BananaPi it is recommended to use the smaller ResNet with Int8-TFLite optimization and for coral boards Int8-TFLite optimization should

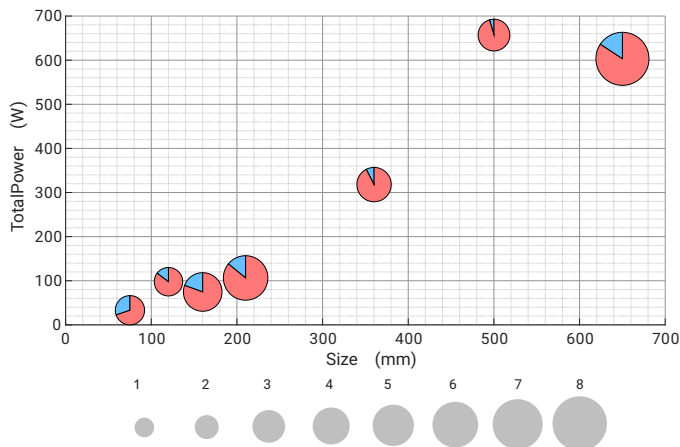


Figure 12. Total Power vs. Quadrotor Size at hover. Each sample is a pie chart which shows the percentage of power consumed by the motors in red and compute and sensing power in blue. The radius of the pie chart is proportional to the power efficiency (in g/W and is given as the ratio of hover thrust to hover power). Refer to the legend on the bottom (gray circles) with the numbers on top indicating power efficiency in g/W.

be used. For the medium sized board TX2 it is advisable to use the smaller SqueezeNet without any optimization (however, we did not explore TensorRT optimization which can result in a 5-20 \times speedup since we limit our analysis to the TensorFlow on PythonTM only environment for ease of use). Again, for larger computers such as a laptop (expect similar performance with a NUC + Jetson Xavier AGX) and a desktop, the smaller SqueezeNet model gives the best performance.

To gather insight into which network dimension (width versus depth) affects the speed the most and to observe the trend on different computers, we measured the speed by varying depth only (Fig. 11a), width only (Fig. 11b) and width + depth together (Fig. 11c). All the networks used in this experiment were VanillaNet₁ (PS \times 1) for consistency. One would expect that increasing depth and width should result in a drop in speed, but this is not always the case. The speeds peak for specific depth/width values (different for different computers) when the perfect balance of memory accesses, size of convolution filters and OPs is achieved. This is more prominent in smaller computers as compared to larger ones (laptop and desktop). This has to be carefully considered when designing networks for use on aerial robots. Also, note that the rate of drop in FPS is more significant with increase in width. A similar trend is observed in Fig. 11c. When the depth is increased and the width is decreased there is a sudden increase in FPS towards the end. Finally, performance using Coral with NanoPi can give higher FPS compared to the desktop or laptop when the models are smaller. Hence, using multiple NanoPi + Coral configurations can be more efficient (in-terms of weight and volume) on larger aerial robots as well.

We categorized quadrotors into six standard configurations based on their size – from 75 mm to 650 mm (pico to large sized), abiding the SWAP constraints. Refer to Table IX and Fig. 12. Each quadrotor is configured with a suitable motor, propeller size and most powerful computer that fits the

respective quadrotor frame. We define ρ as the ratio of hover power with and without computer. Most literature only talks about the amount of power the on-board computer uses, but this only highlights half of the story. Adding a computer to a quadrotor (to make it autonomous) not only consumes power from the battery but also adds weight which in-turn makes the motor power consumption higher at hover (one would need more thrust to keep the quadrotor flying). Generally, the power efficiency (thrust to power ratio in g/W) decreases with an increase in thrust aggravating the situation further. The essence of this is quantified by ρ which indicates how much more power the autonomous quadrotor uses compared to a manual one with the same configuration (excluding computer and sensors). A value of $\rho = 1.0$ is the theoretical best, and the larger the value (above 1.0), the more inefficient the setup. One could clearly observe that the amount of thrust directly increases with quadrotor size as it should but the thrust-to-weight ratio follows a parabolic curve (opening at the top) which achieves a maximum value at 210 mm size. This is due to efficient motor design perfected for racing quadrotors. We also observe that for smaller quadrotors (≤ 75 mm) the power overhead due to adding the computer is significant (as high as 92%). The power overhead decreases and then increases again as size increases due to the addition of multiple computers at 650 mm size. Also, note that for aggressively flying quadrotors the value of ρ will decrease significantly since we choose the most efficient hovering motors available on the market to maximize battery life. Fig. 1 shows four different sized quadrotors, computers and other commonly used quadrotor electronics for a size comparison. We also show how small a hardware designed from the ground-up (Snapdragon flight) can be. We exclude this from our discussions since the smaller model of Snapdragon flight is currently phased out by Qualcomm. Finally, we also experimented with the Sipeed Maix Bit which is a low power neural network accelerator (< 1 W of power) weighing only 20 g with a camera and which can be used on smaller sized quadrotors. However, due to lack of ease of use we exclude it from our discussion.

C. Trajectory Evaluation

The predictions \tilde{h} (VanillaNet₄ T \times 2, S \times 2 large model trained using supervised l_2 loss) are obtained every four frames and are integrated using dead-reckoning to obtain the final trajectory. The trajectory is aligned with the ground truth and evaluated using the approach given in [61]. The errors in individual axes (l_1 distance) and all axes (RMSE) are given for various trajectories in Table X and are illustrated in Fig. 13. We notice that even with simple dead-reckoning we obtain an RMSE of less than 3% of the trajectory length highlighting the robustness of the proposed PRGFlow.

VI. SUMMARY AND DIRECTIONS FOR FUTURE WORK

A summary of our observations are given below:

- The effect of the choice of network architecture is fairly significant on the accuracy, number of parameters and

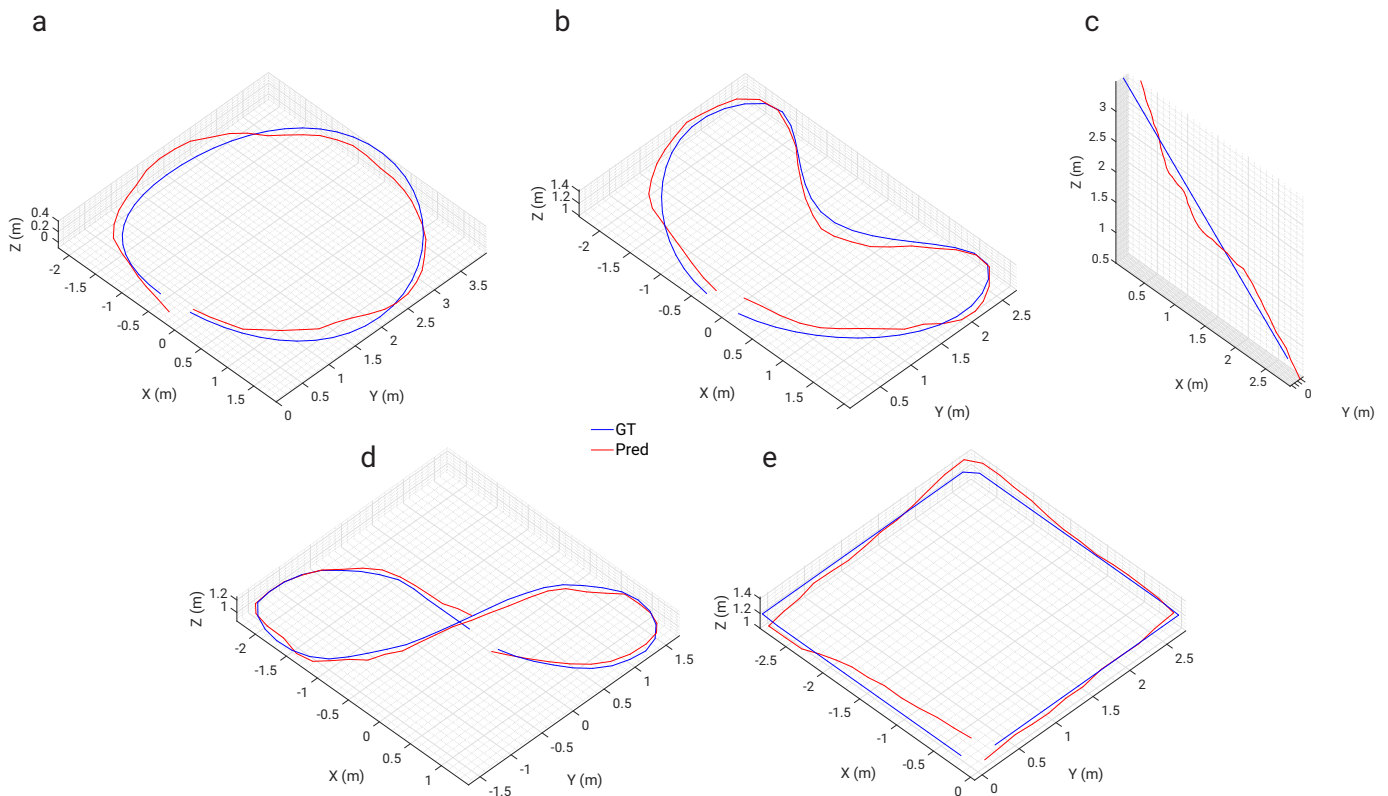


Figure 13. Comparison of trajectory obtained by dead-reckoning (red) our estimates with respect to the ground truth (blue) for quadrotor flight in various trajectory shapes. (a) Circle, (b) Moon, (c) Line, (d) Figure8 and (e) Square.

OPs and needs to be carefully considered when designing a network for deployment on an aerial robot.

- Although number of parameters is inversely correlated to speed, sometimes increasing the number of parameters (depth or width) can achieve a speedup due to a better balance in memory accesses, size of convolution filters and OPs which is more significant in smaller computers than larger ones.
- It is advisable to use supervised methods for odometry estimation since they are much easier to train and are generally more accurate than their unsupervised counterparts. The loss in accuracy from simulation training to real world is insignificant if enough variation in samples is used which have real world images.
- For odometry related regression problems, the simplest method of dropping weights works as well as more complex distillation methods.
- Accelerated instruction sets can provide huge speedups and should not be neglected for neural network based methods.
- Lastly, deep learning based odometry have two main advantages: they are generally faster due to hardware parallelization and fail more gracefully (work better in adverse scenarios) when compared to their classical counterparts.
- A combination of both classical and deep learning methods to solve a problem generally would yield better explainability and performance gains.

Based on the observations and empirical analyses we hope the following directions for future work can further enhance ego-motion/odometry capabilities using deep learning.

- Crafting better supervised loss functions (probably on manifolds) should yield more accurate and robust results and probably will surpass advantages of un-supervised/self-supervised methods given enough domain variation.
- Formulating problems as nested loops can further simplify problems and provide better explainability than end-to-end deep learning. Also, this enables solving certain problems directly with simple mathematical formulations which can further decrease the number of neurons/OPs required, thereby reducing latency in most cases.
- Utilizing multi-frame constraints can significantly boost performance of deep odometry systems.
- Better architectures need to be designed to fully utilize the capabilities of un-supervised/self-supervised methods of deep odometry estimation.

VII. CONCLUSIONS

We presented a simple way to estimate ego-motion/odometry on an aerial robot using deep learning combining commonly found sensors on-board: a up/down-facing camera, an altimeter source and IMU. By utilizing simple filtering methods to estimate attitude one can obtain “cheap” odometry using attitude compensated

frames as the input to a network. We further provided a comprehensive analysis of warping combinations, network architectures and loss functions. All our approaches were benchmarked on different commonly used hardware with different SWAP constraints for speed and accuracy which we hope will serve as a reference manual for researchers and practitioners alike. We also show extensive real-flight odometry results highlighting the robustness of the approach without any fine-tuning or re-training. Finally, as a parting thought, utilizing deep learning when failure is often expected would most likely lead to more robust system.

ACKNOWLEDGEMENT

The support of the Brin Family Foundation, the Northrop Grumman Mission Systems University Research Program, ONR under grant award N00014-17-1-2622, and the support of the National Science Foundation under grant BCS 1824198 are gratefully acknowledged. We would also like to thank NVIDIA for the grant of a Titan-Xp GPU used in the experiments.

REFERENCES

- [1] Nitin J. Sanket et al. GapFlyt: Active vision based minimalist structure-less gap detection for quadrotor flight. *IEEE Robotics and Automation Letters*, 3(4):2799–2806, Oct 2018.
- [2] Teodor Tomic et al. Toward a fully autonomous uav: Research platform for indoor and outdoor urban search and rescue. *IEEE robotics & automation magazine*, 19(3):46–56, 2012.
- [3] Nathan Michael et al. Collaborative mapping of an earthquake-damaged building via ground and aerial robots. *Journal of Field Robotics*, 29(5):832–841, 2012.
- [4] Cornelia Fermüller. Navigational preliminaries. In Y. Aloimonos, editor, *Active Perception*. Lawrence Erlbaum Associates, 1993.
- [5] KN McGuire, C De Wagter, K Tuyls, HJ Kappen, and GCHE de Croon. Minimal navigation solution for a swarm of tiny flying robots to explore an unknown environment. *Science Robotics*, 4(35):eaaw9710, 2019.
- [6] Ji Zhang and Sanjiv Singh. Loam: Lidar odometry and mapping in real-time. In *Robotics: Science and Systems*, volume 2, 2014.
- [7] Ji Zhang and Sanjiv Singh. Visual-lidar odometry and mapping: Low-drift, robust, and fast. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2174–2181. IEEE, 2015.
- [8] Fabio Morbidi, Randy A Freeman, and Kevin M Lynch. Estimation and control of uav swarms for distributed monitoring tasks. In *Proceedings of the 2011 American Control Conference*, pages 1069–1075. IEEE, 2011.
- [9] Aaron Weinstein, A Cho, Giuseppe Loianno, and Vijay Kumar. Vio-swarm:: A swarm of 250g quadrotors. *IEEE RA-L Robotics and Automation Letters*, 2018.
- [10] Daigo Shishika and Derek A Paley. Mosquito-inspired distributed swarming and pursuit for cooperative defense against fast intruders. *Autonomous Robots*, 43(7):1781–1799, 2019.
- [11] Morgan Quigley, Kartik Mohta, Shreyas S Shivakumar, Michael Watterson, Yash Mulgaonkar, Mikael Arguedas, Ke Sun, Sikang Liu, Bernd Pfrommer, Vijay Kumar, et al. The open vision computer: An integrated sensing and compute system for mobile robots. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 1834–1840. IEEE, 2019.
- [12] Anastasios I Mourikis and Stergios I Roumeliotis. A multi-state constraint kalman filter for vision-aided inertial navigation. In *Proceedings 2007 IEEE International Conference on Robotics and Automation*, pages 3565–3572. IEEE, 2007.
- [13] Nitin J. Sanket, Chethan M. Parameshwara, Chahat Deep Singh, Ashwin V. Kuruttukulam, Cornelia Fermüller, Davide Scaramuzza, and Yiannis Aloimonos. Evdodgenet: Deep dynamic obstacle dodging with event cameras, 2019.
- [14] Dominik Honegger, Lorenz Meier, Petri Tanskanen, and Marc Pollefeys. An open source and open hardware embedded metric optical flow cmos camera for indoor and outdoor applications. In *2013 IEEE International Conference on Robotics and Automation*, pages 1736–1741. IEEE, 2013.
- [15] Takafumi Taketomi et al. Visual slam algorithms: a survey from 2010 to 2016. *IPSN Transactions on Computer Vision and Applications*, 9(1):16, June 2017.
- [16] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *IEEE transactions on robotics*, 31(5):1147–1163, 2015.
- [17] Jakob Engel, Thomas Schöps, and Daniel Cremers. Lsd-slam: Large-scale direct monocular slam. In *European conference on computer vision*, pages 834–849. Springer, 2014.
- [18] Sean L Bowman, Nikolay Atanasov, Kostas Daniilidis, and George J Pappas. Probabilistic data association for semantic slam. In *2017 IEEE international conference on robotics and automation (ICRA)*, pages 1722–1729. IEEE, 2017.
- [19] Stefan Leutenegger, Simon Lynen, Michael Bosse, Roland Siegwart, and Paul Furgale. Keyframe-based visual-inertial odometry using nonlinear optimization. *The International Journal of Robotics Research*, 34(3):314–334, 2015.
- [20] Michael Bloesch et al. Robust visual inertial odometry using a direct ekf-based approach. In *2015 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 298–304. IEEE, 2015.
- [21] Jakob Engel et al. Direct sparse odometry. *IEEE transactions on pattern analysis and machine intelligence*, 2017.
- [22] Tong Qin et al. VINS-Mono: A robust and versatile monocular visual-inertial state estimator. *IEEE Transactions on Robotics*, 34(4):1004–1020, 2018.
- [23] Huai-Jen Liang et al. SalientDSO: Bringing attention to direct sparse odometry. *IEEE Transactions on Automation Science and Engineering*, 2019.
- [24] Alex Kendall, Matthew Grimes, and Roberto Cipolla. Posenet: A convolutional network for real-time 6-dof camera relocalization. In *Proceedings of the IEEE international conference on computer vision*, pages 2938–2946, 2015.
- [25] Alex Kendall and Roberto Cipolla. Geometric loss functions for camera pose regression with deep learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5974–5983, 2017.
- [26] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1851–1858, 2017.
- [27] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1983–1992, 2018.
- [28] Nan Yang, Lukas von Stumberg, Rui Wang, and Daniel Cremers. D3vo: Deep depth, deep pose and deep uncertainty for monocular visual odometry, 2020.
- [29] Ronald Clark, Sen Wang, Hongkai Wen, Andrew Markham, and Niki Trigoni. Vinet: Visual-inertial odometry as a sequence-to-sequence learning problem. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [30] Liming Han, Yimin Lin, Guoguang Du, and Shiguo Lian. Deepvio: Self-supervised deep learning of monocular visual inertial odometry using 3d geometric constraints, 2019.
- [31] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013.
- [32] Michael Burri, Janosch Nikolic, Pascal Gohl, Thomas Schneider, Joern Rehder, Sammy Omari, Markus W Achtelik, and Roland Siegwart. The euroc micro aerial vehicle datasets. *The International Journal of Robotics Research*, 35(10):1157–1163, 2016.
- [33] Jakob Engel, Vladyslav Usenko, and Daniel Cremers. A photometrically calibrated benchmark for monocular visual odometry. *arXiv preprint arXiv:1607.02555*, 2016.
- [34] Bernd Pfrommer, Nitin Sanket, Kostas Daniilidis, and Jonas Cleveland. Pcnosyvio: A challenging visual inertial odometry benchmark. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3847–3854. IEEE, 2017.
- [35] Jeffrey Delmerico and Davide Scaramuzza. A benchmark comparison of monocular visual-inertial odometry algorithms for flying robots. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2502–2509. IEEE, 2018.
- [36] Simone Bianco, Remi Cadene, Luigi Celona, and Paolo Napoletano. Benchmark analysis of representative deep neural network architectures. *IEEE Access*, 6:64270–64277, 2018.

- [37] Alfredo Canziani, Adam Paszke, and Eugenio Culurciello. An analysis of deep neural network models for practical applications. *arXiv preprint arXiv:1605.07678*, 2016.
- [38] Kirk YW Schepers and Guido CHE de Croon. Evolution of robust high speed optical-flow-based landing for autonomous mavs. *Robotics and Autonomous Systems*, 124:103380, 2020.
- [39] Chen-Hsuan Lin and Simon Lucey. Inverse compositional spatial transformer networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [40] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [41] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [42] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [43] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 116–131, 2018.
- [44] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [45] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.
- [46] Hang Zhao, Orazio Gallo, Iuri Frosio, and Jan Kautz. Loss functions for image restoration with neural networks. *IEEE Transactions on computational imaging*, 3(1):47–57, 2016.
- [47] D. Sun et al. A quantitative analysis of current practices in optical flow estimation and the principles behind them. *International Journal of Computer Vision*, 106(2):115–137, 2014.
- [48] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [49] J. Barron. A general and adaptive robust loss function. *CVPR*, 2019.
- [50] Sebastian OH Madgwick, Andrew JL Harrison, and Ravi Vaidyanathan. Estimation of imu and mag orientation using a gradient descent algorithm. In *2011 IEEE international conference on rehabilitation robotics*, pages 1–7. IEEE, 2011.
- [51] Sujith Ravi. Projectionnet: Learning efficient on-device deep networks using neural projections. *arXiv preprint arXiv:1708.00630*, 2017.
- [52] Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? In *Advances in neural information processing systems*, pages 2654–2662, 2014.
- [53] B Srinivasa Reddy and Biswanath N Chatterji. An fft-based technique for translation, rotation, and scale-invariant image registration. *IEEE transactions on image processing*, 5(8):1266–1271, 1996.
- [54] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. Surf: Speeded up robust features. In *European conference on computer vision*, pages 404–417. Springer, 2006.
- [55] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. Orb: An efficient alternative to sift or surf. In *2011 International conference on computer vision*, pages 2564–2571. Ieee, 2011.
- [56] Deepak Geetha Viswanathan. Features from accelerated segment test (fast). *Homepages. Inf. Ed. Ac. Uk*, 2009.
- [57] Stefan Leutenegger, Margarita Chli, and Roland Y Siegwart. Brisk: Binary robust invariant scalable keypoints. In *2011 International conference on computer vision*, pages 2548–2555. Ieee, 2011.
- [58] Christopher G Harris, Mike Stephens, et al. A combined corner and edge detector. In *Alvey vision conference*, volume 15, pages 10–5244. Citeseer, 1988.
- [59] Matthias Faessler, Davide Falanga, and Davide Scaramuzza. Thrust mixing, saturation, and body-rate control for accurate aggressive quadrotor flight. *IEEE Robotics and Automation Letters*, 2(2):476–482, 2016.
- [60] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 224–236, 2018.
- [61] Zichao Zhang and Davide Scaramuzza. A tutorial on quantitative trajectory evaluation for visual (-inertial) odometry. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7244–7251. IEEE, 2018.