

Secure Deep Learning-based Distributed Intelligence on Pocket-sized Drones

Elia Cereda
IDSIA, USI-SUPSI, Switzerland
elia.cereda@idsia.ch

Alessandro Giusti
IDSIA, USI-SUPSI, Switzerland
alessandro.giusti@idsia.ch

Daniele Palossi
IDSIA, USI-SUPSI, Switzerland
IIS, ETH Zurich, Switzerland
daniele.palossi@idsia.ch

Abstract

Palm-sized nano-drones are an appealing class of edge nodes, but their limited computational resources prevent running large deep-learning models onboard. Adopting an edge-fog computational paradigm, we can offload part of the computation to the fog; however, this poses security concerns if the fog node, or the communication link, can not be trusted. To tackle this concern, we propose a novel distributed edge-fog execution scheme that validates fog computation by redundantly executing a random subnetwork aboard our nano-drone. Compared to a State-of-the-Art visual pose estimation network that entirely runs onboard, a larger network executed in a distributed way improves the R^2 score by +0.19; in case of attack, our approach detects it within 2 s with 95% probability.

Supplementary material

In-field experiments and system demonstration video:
<https://youtu.be/QwTiigAs4cA>.

1 Introduction

With an increasing number of Internet-of-Things-capable (IoT) end devices, from tiny headphones to autonomous cars, the *edge-fog paradigm* permeates almost any civil and industrial application [6]. This paradigm exploits distributed computation with a resource-limited edge device close to the source of data, combined with a more capable remote fog node able to overcome the memory and computational limits of the former node. As the edge node, this work considers a novel and appealing cyber-physical system: a pocket-sized quadrotor or nano-drone. Thanks to their sub-10 cm diameter, nano-drones have the potential to unlock unprecedented application scenarios, such as the exploration of narrow, cluttered or GPS-denied environments, safe human-robot interaction, and intelligent ubiquitous IoT sensing. Additionally,

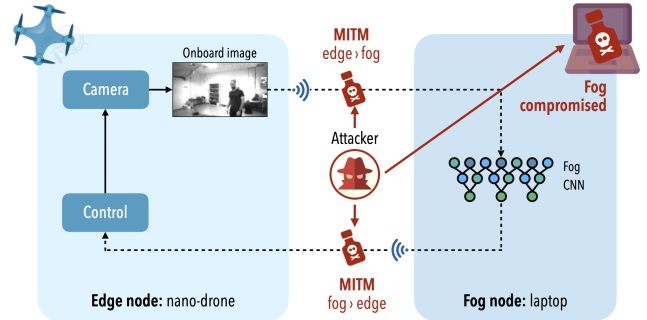


Figure 1. The threat model of our edge-fog use case.

these platforms are much less expensive than larger ones due to simplified electronics and mechanics.

However, with their small form factor comes their main limitation: they can only host ultra-constrained processors and sensors, i.e., simple ultra-low power microcontroller units (MCUs) with a sub-100 mW power envelope. This limitation poses two challenges; on the one hand, they have extremely limited computational power and memory to run aboard complex real-time algorithms. To partially overcome this issue, nano-drones often adopt navigation engines based on small convolutional neural networks (CNNs) [12, 7], due to their compelling trade-off between accuracy and computational cost, compared to geometrically precise computer vision-based [10] and predictive methods [9].

On the other hand, the security features nano-drones can offer are still quite rudimentary and rarely addressed by State-of-the-Art (SoA) solutions [18]. An edge-fog scheme introduces two significant vulnerabilities in the remote fog node and the communication channel: compromising either of them, e.g., with fog malware, man-in-the-middle, data infiltration, or data-in-transit attacks, allows an attacker to send malicious commands to the edge and take control of its behavior. Our work tries to mitigate both vulnerabilities.

We address the human pose estimation task and embody the edge-fog paradigm with a nano-drone (edge) and a WiFi-connected commodity laptop (fog), which acts as a powerful remote brain to boost the perception capabilities of the edge node, as depicted in Figure 1. In our case, the edge-fog execution scheme enables execution of more complex CNNs, otherwise unaffordable aboard the nano-drone within its real-time constraints. We design a three-stage distributed

MobileNetV2-based [14] CNN to perform our vision-based task. The first compression stage transforms a large input image into a smaller tensor; a central backbone performs the most computationally-intense processing; a final reduction stage produces the output. Then, we design the central backbone with a multi-branch structure, where multiple tensors are computed independently and in parallel before being reduced in the final stage. We execute the first and last stages on the nano-drone while we offload the heavy multi-branch backbone to the fog node.

Our strategy enables an additional level of security between edge and fog nodes, which can be stacked on top of traditional encryption mechanisms if affordable by the edge MCU. While the fog executes all branches for every input image, the edge also executes one random branch at a time and checks its resulting tensor against the one received from the fog, embodying a probabilistic security mechanism. If the check is successful, the edge uses all the received tensors (one per branch) to run the final reduction. Otherwise, it falls back to a small field-proofed CNN able to run in real-time aboard a nano-drone, i.e., the PULP-Frontnet [12]. The PULP-Frontnet acts as a backup solution, taking control when either the information received by the fog is compromised or when the communication channel is disrupted/unavailable (e.g., jamming). This shallow CNN trades the regression accuracy w.r.t. the more accurate remote multi-branch CNN for lower complexity, making it our “last resort” to complete or gracefully abort the mission.

We demonstrate our general approach by showing the perceptible benefit of a multi-branch MobileNetV2-based CNN, which achieves a mean $+0.19 R^2$ score increase compared to the PULP-Frontnet baseline. Then, we introduce a simple model accounting for computation and communication overheads to identify the optimal points to split our CNN between compression, backbone, and reduction stages to maximize the overall throughput up to ~ 8 frame/s. Finally, we deploy our distributed system in the field, assessing its security functionalities (see in-field video) and closed-loop performance with a 31.3% (horizontal) and 34.6% (angular) reduction of the control error compared to the baseline.

2 Related work

Autonomous navigation algorithms executed aboard nano-robotics platforms are subject to severe computational and memory constraints. SoA CNNs deployed on nano-drones exploit novel parallel ultra-low-power systems-on-chip (SoCs) [1, 11, 12], but they can only execute much simpler real-time workloads compared to those running on larger-scale UAVs. Therefore, limiting the complexity of the achievable tasks and the quality of the obtained results. Among these works, PULP-Frontnet [12] is a monocular CNN for the human pose estimation task with a computational cost of 14.7 MMAC (millions of multiply-accumulate operations). It achieves an inference rate of 48 frame/s while consuming 96 mW, aboard a Crazyflie 2.1 nano-drone.

By comparison, common visual CNN architectures [3, 13, 5, 20] deployed on larger drones have orders of magnitude higher computation workloads. For example, Foehn et al. [3] tackle the task of gate pose estimation in an autonomous

drone race using a 1930 MMAC per inference CNN running on an Nvidia Jetson Xavier board (~ 20 W). These systems either depend on more powerful on-board processors than available on nano-drones [3, 5] or offload the computation to a remote node and thus suffer from the security vulnerabilities we address in this work [13, 20].

Several ad-hoc neural network architectures have been proposed for efficient execution while retaining SoA task performance, such as MobileNetV2 [14] or EfficientNetV2 [15]. We select MobileNetV2, i.e., one of the most widespread CNN for embedded devices, as the base architecture for our work. Given its design for smartphone-class devices with 90 MMAC per inference, we leverage a fog node to achieve real-time throughput and SoA regression performance on our task. Additionally, despite the specific architecture used in our use case, we present a general approach that can be applied to other CNNs and for different tasks.

We extend MobileNetV2 with a distributed execution scheme based on multiple independent branches to achieve the desired security properties. Multi-branch architectures have been adopted in several recent works for various reasons, such as *i*) reducing the computational cost by dynamically executing only subsets of the branches [8, 16], *ii*) to estimate uncertainty in network predictions [17]; *iii*) to prevent adversarial attacks [4]. Crucially, these works all focus on models running on a single device, while our work is the first to employ a multi-branch architecture to detect attacks on a distributed edge-fog system, where an agile nano-drone embodies the edge node.

3 System design

3.1 Use case

In the envisioned scenario, the nano-drone (edge node) is tasked with a vision-based perception task, i.e., human pose estimation. A more capable fog node, represented by a resource-unconstrained remote laptop, provides additional computational power to the nano-drone to increase its execution performance. Edge and fog nodes share a bi-directional communication channel, i.e., WiFi, on which they can continuously stream data. We select two different deep-learning models to perform the pose estimation task. The former is a lightweight (300 k parameters) convolutional neural network (CNN) called PULP-Frontnet [12] and previously demonstrated running in real-time aboard nano-drones. The latter is a more memory-/computationally-demanding (1.8 M parameters) but also more accurate CNN, based on the MobileNetV2 model [14]. The MobileNetV2-based workload is distributed between edge and fog, where most of it is executed by the more capable remote fog, leaving a smaller fraction of computation on the memory/processor-limited edge. The lightweight PULP-Frontnet is executed entirely on the edge, only if the communication channel or the remote computation becomes compromised, as detailed in the following.

3.2 CNN design

The traditional execution scheme of many vision-based CNNs can be mapped into three main computational stages, as shown in Figure 2. *i*) an initial compression stage, where a high-resolution input image is compressed to a smaller-size tensor; *ii*) the execution of a central computationally-

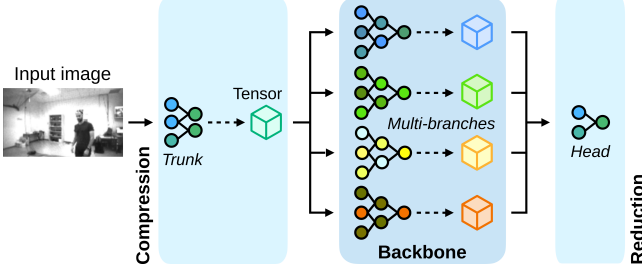


Figure 2. Vision-based CNN’s three main stages.

demanding backbone; and *iii*) a final reduction stage, e.g., a fully connected layer, producing the final output. Under this consideration, we design our MobileNetV2-based CNN as a multi-branch model, where the central backbone is split into \mathcal{N} independent branches. Each branch is fed with the same input tensor, i.e., the output of the initial compression stage – called *trunk* – and produces its output tensor. The last reduction stage, called *head*, takes as input the concatenation of all tensors resulting from the multi-branch backbone and produces the final output for the robot controller.

We apply this general approach to the task of human pose estimation. Our CNN takes as input one gray-scale 160×96 px image and estimates the pose of the human subject relative to the drone. The pose components are represented as four independent regression outputs, which correspond to the three Cartesian coordinates of the subject’s position in 3D space (x, y, z) w.r.t. the drone’s horizontal frame and the subject’s relative orientation w.r.t. the drone’s yaw, ϕ . Finally, we employ this relative pose (x, y, z, ϕ) to perform the “follow-me” application aboard our nano-drone.

3.3 Threat model

The success of our mission depends on *i*) correct execution of both edge and fog computation, and *ii*) reliable and unaltered data exchange between edge and fog nodes. We focus on the two vulnerabilities shown in Figure 1, in which a malicious actor either compromises the entire fog node or the communication channel, e.g., man-in-the-middle, data infiltration, and data-in-transit hacked attacks. Either scenario would allow the attacker to poison the information transmitted by the fog node and thus take control of the nano-drone. To mitigate this threat, we introduce a probabilistic security mechanism that relies on redundant execution of a small amount of computation on the edge node to assess the trustworthiness of the data received from the fog.

Our approach is orthogonal to computationally-demanding traditional countermeasures, such as communication encryption (not always affordable on MCU-limited edge nodes), providing a *defense-in-depth* strategy. Finally, the last attack we consider is complete loss of the edge-fog communication channel (e.g., radio jamming). In this case, we resort to safe execution of the PULP-Frontnet CNN entirely aboard our nano-drone, degrading pose estimation accuracy but without jeopardizing the mission.

3.4 Proposed strategy

We propose a strategy based on partial computation redundancy to enable an onboard probabilistic security mech-

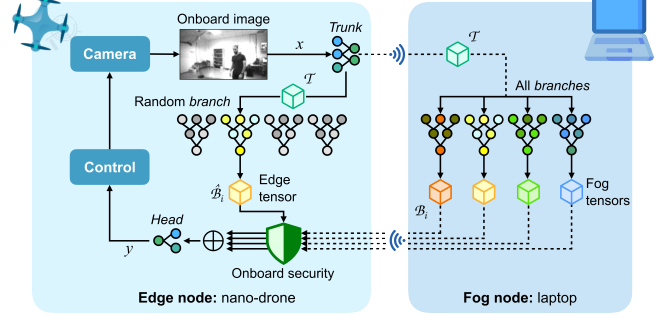


Figure 3. Our distributed probabilistic security strategy.

anism while increasing the edge inference throughput thanks to the computationally-capable fog. Figure 3 shows the proposed system. The execution starts from the edge node running the first compression stage (*trunk* CNN) on the input image (x) and producing a compressed tensor \mathcal{T} to reduce the streaming data between edge and fog nodes:

$$\mathcal{T} = \text{trunk}(x) \quad (1)$$

\mathcal{T} is then forwarded via WiFi to the fog and used locally on the edge. The fog node executes the entire multi-branch backbone for each input tensor \mathcal{T} . The CNN running on the fog is the backbone (central part) of a MobileNetV2 CNN composed of \mathcal{N} independent branches, where \mathcal{N} is a hyper-parameter of the model. Each branch receives the same input \mathcal{T} and produces one unique feature tensor \mathcal{B}_i :

$$\mathcal{B}_i = \text{branch}_i(\mathcal{T}) \quad \text{for } i = 1 \dots \mathcal{N} \quad (2)$$

Feature tensors from all branches are then forwarded via WiFi to the edge.

On the edge node, only one random branch j is executed per input image to minimize the computation burden. A cryptographic random number generator ensures that an attacker cannot predict the sequence of verified branches. The selected branch is fed with \mathcal{T} and produces $\hat{\mathcal{B}}_j$, which is compared against the fog’s corresponding computation \mathcal{B}_j . If \mathcal{B}_j is bit-by-bit identical to $\hat{\mathcal{B}}_j$, then all the tensors received from the fog are considered trustable and concatenated (\oplus operator) before being reduced in the final stage:

$$y = \text{head}\left(\bigoplus_i \mathcal{B}_i\right) \quad (3)$$

where the *head* CNN computes the final output y on the edge.

Under this system design, the attacker maximizes their chances of avoiding detection by tampering with only a single feature tensor \mathcal{B}_k . This translates to a probability $1/\mathcal{N}$ of detecting an attack in any given frame (i.e., the probability that the edge node executes the same branch corrupted by the attacker $j = k$). Nevertheless, our detection scheme is so lightweight that it can run continuously as part of a robot’s closed-loop controller. Therefore, the attacker must continue tampering over time to keep control of the nano-drone, but our attack detection probability converges to 1.0 as more frames are verified.

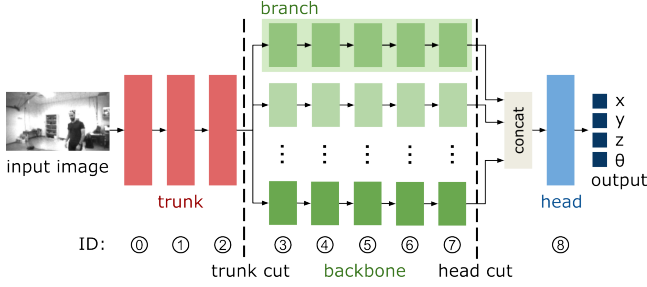


Figure 4. Multi-branch MobileNetV2 CNN architecture.

3.5 Deployment

Robotic platform. Our edge node is embodied by a commercial off-the-shelf (COTS) Crazyflie 2.1 nano-quadrotor, an open-source 27 g palm-sized nano-drone produced by Bitcraze. An STM32 MCU performs low-level flight control, while an AI-deck COTS expansion board provides high-level vision-based perception with an additional GWT GAP8 SoC. The GAP8 is a multi-core processor featuring 9 RISC-V-based cores, split between an eight-core *cluster* optimized for parallel computation of compute-intensive workloads and a single-core *fabric controller* that manages interfaces with external peripherals and on-chip memories. The on-chip memory hierarchy is organized into two levels: 64 kB low-latency L1 memory and 512 kB L2 memory. This expansion board also provides off-chip memories (8 MB DRAM and 64 MB Flash), a monochrome QVGA Himax HM01B0 camera, and an ESP32-based Ublox NINA-W102 WiFi module. The lack of hardware floating-point units and data cache memories on the GAP8 SoC dictate, respectively, integer-quantized arithmetic and explicit data management between memories. We leverage an ad-hoc CNN deployment pipeline to generate C code that addresses these concerns [12]. The fog node also runs its workload in fixed-point arithmetic to simplify bit-by-bit tensor comparisons.

Multi-branch architecture.

We design our CNN starting from the widely adopted MobileNetV2 [14], which requires 89 MMAC operations per inference. First, we define a coarse split of the three CNN segments we introduced in Section 3.2: the trunk, a multi-branch backbone, and the final head, see Figure 4. Branching the execution from the beginning of the CNN, which means discarding the trunk, would lead to poor regression performance and a waste of resources by neglecting an initial trainable path shared by all branches [19]. Therefore, initially, we define as the *cutting point* between the trunk and backbone the first layer having an output tensor smaller than the input image (~ 15 kB), which results in a trunk of 25.7 MMAC operations. The remaining 63.3 MMAC are equally distributed among all \mathcal{N} branches, splitting the number of output channels. For example, in the case of $\mathcal{N} = 8$, a layer producing an output tensor of 32 channels in the original CNN is split into 8 layers, each outputting a tensor of 4 channels.

Figure 5 shows four configurations in the number of backbone’s branches, i.e., 1, 2, 4, and 8. For each configuration, we report on the primary y-axis the probability of detecting an attack within a dt time of 0.5, 1, 1.5, and 2 s; the higher the

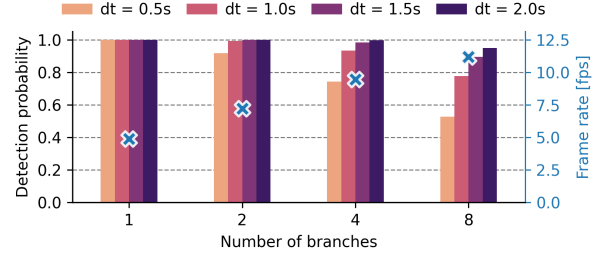


Figure 5. Detection probability vs. inference frame-rate.

dt , the higher the probability of converging to 1. On the secondary y-axis, we show the achievable frame rate for each configuration; more branches result in less computation on the edge node and lower detection probability. The real-time constraint posed by our cyber-physical system suggests selecting a configuration with a frame rate as higher as possible, with a lower bound of ≈ 10 frame/s [12]. Therefore, we use the configuration $\mathcal{N} = 8$ for our experiments, which also achieves a detection probability of 95% with a $dt=2$ s.

Finally, in Figure 6, we thoroughly analyze how different cutting points affect the edge execution time and TX/RX communication latencies between edge and fog, for which we assume a negligible execution time having multiple orders of magnitude more capable compute unit. Layers in the MobileNetV2 architecture are grouped in *inverted residual blocks*, inside which tensor sizes follow an expansion-compression pattern. For this reason, we force cutting points to be placed at block boundaries, which we show in Figure 4 with 9 IDs. Earlier cutting points in the trunk, y-axis in Figure 6, minimize its execution time (A) but inflates the TX communication latency (B) as more data are streamed to the fog. Similarly, later cutting points of the head, x-axis in Figure 6, increase the execution time of the trunk (A), the backbone’s one branch (C), but reduce the latencies for both RX communication (D) and the head execution (E). End-to-end latency (F) considers that the edge node executes the backbone branch in parallel with communication:

$$\text{end-to-end} = \text{trunk} + \max(\text{backbone}, \text{TX} + \text{RX}) + \text{head} \quad (4)$$

The resulting optimal configuration is trunk cutting point ID = 3 and head cutting point ID = 8. Finally, we deploy and profile this optimal model on the GAP8 SoC, obtaining a $\pm 5\%$ execution time compared to the estimates in Figure 6.

4 Results

We report three experiments. The first shows the improved regression performance of our proposed multi-branch model. The second shows the resulting improvement of in-field robot behavior. The third demonstrates the attack detection capabilities of the proposed security scheme.

4.1 Regression performance

We use the dataset presented in [2] for our training and regression performance analysis. This dataset contains images acquired with the same robotic platform employed in our work and labels of the absolute pose of human subjects and a nano-drone. Data are collected from two indoor laboratories equipped with a motion capture system (mocap),

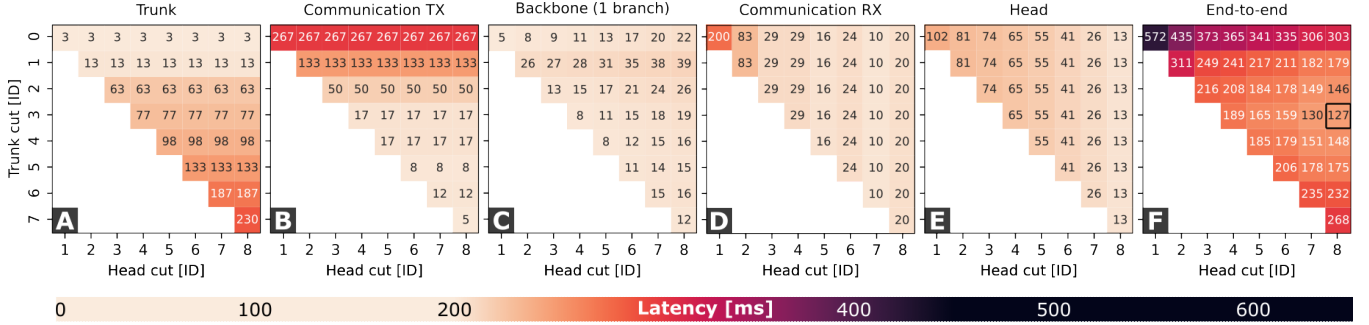


Figure 6. Latency induced by each component, for every combination of cut points.

accounting for 12k images collected with 17 distinct human subjects (age, height, etc.). Three subjects (4.7k samples) form our test set, while the remaining 14 subjects (7.3k) are used for training (90%) and validation (10%). Starting from models pre-trained on ImageNet, we train for 100 epochs with the Adam optimizer and learning rate 10^{-3} to minimize the L1 loss of the relative pose.

Regression performance represents the ability of the model to estimate the subject’s relative pose accurately. We quantify it through the R^2 score (or *coefficient of determination*) computed for each output variable. R^2 is a standard adimensional metric representing the fraction of variance in the target variable explained by the model. An ideal model that perfectly estimates the variable yields $R^2 = 1$; a trivial model that always returns the test set average yields $R^2 = 0.0$.

Figure 7 compares the proposed fog+edge multi-branch models against the SoA edge-only PULP-Frontnet and the fog-only original MobileNetV2. Each model is trained five times with randomly-initialized parameters. The original MobileNetV2 shows a $\sim 7\times$ increase in network parameters and computation workload, which pays off with an increase of more than +0.20 in median R^2 score compared to PULP-Frontnet, justifying the use of a fog node and larger models. Compared to the original MobileNetV2, multi-branch models yield similar (or slightly improved) computation workloads and R^2 scores on all output variables. Thus, introducing the multi-branch architecture enables our security application without impacting other aspects of the system. Figure 8 compares individual predictions of the 8-branch MobileNetV2 against PULP-Frontnet, showing a significantly lower regression error on all four pose components.

4.2 Control performance

This experiment evaluates our model’s tracking accuracy when used in a fully-autonomous in-field closed-loop control system. We reproduce the testing setup of our baseline [12], where the human subject, never seen by the models during training, moves along a predefined path of increasing difficulty. At the same time, the autonomous drone is tasked to stay in front of the subject at a constant distance of 1.5 m, therefore behaving in the “follow-me” application. We compare two models: the SoA PULP-Frontnet baseline, running entirely on the edge node, and the 8-branch MobileNetV2 model, running remotely on the fog node, which is also affected by the end-to-end communication latency. While a mocap system tracks both subject and drone poses, we per-

Table 1. Control error.

	PULP-Frontnet [12]	MobileNetV2 (8 branches)	Mocap
Control error			
e_{xy} [m]	0.99	0.68	0.18
e_θ [rad]	0.75	0.49	0.21

form three flights per model ($3 \times 2 = 6$ flights) plus one additional flight in which the controller relies on the subject’s ground-truth pose (as measured by the mocap), representing the behavior with a perfect predictor.

Table 1 measures the resulting control performance through two error terms: e_{xy} , the drone’s mean horizontal distance from its desired position (i.e., 1.5 m in front of the subject), and e_θ , the drone’s mean absolute angular error from its desired orientation (i.e., looking directly at the subject). The mocap-based flight represents the control error lower bound achievable with perfect predictions. We observe that the 8-branch MobileNetV2 model significantly improves on the baseline, reducing both errors by -30%. The supplementary video also shows these experiments, highlighting the improved control accuracy.

4.3 Onboard security demonstration

In the supplementary video, we demonstrate our system’s security capabilities in an in-field qualitative demonstration. Initially, the nano-drone performs the “follow-me” task in normal conditions without attacks nor security mechanisms in place. This results in good tracking performance, with the nano-drone able to precisely follow the human subject (never seen in training). After a few tens of seconds of flight, we simulate an attack in which the fog node continuously sends malicious tensors to the edge. This malicious tensor encodes a final regression output of $x < 1.5$ m, which forces the drone to fly away from the subject, i.e., perceived as too close. Then, the same attack is repeated, but this time having the probabilistic security mechanism active. In this case, the attack is detected within 1 s, and the edge node reacts by triggering a predefined emergency behavior, i.e., hovering in place. Once the attack terminates, the fog tensor returns to be trusted by the onboard security mechanism, and the nano-drone resumes following the human subject.

5 Conclusion

We present a probabilistic security mechanism built on top of an edge-fog system embodied by a resource-

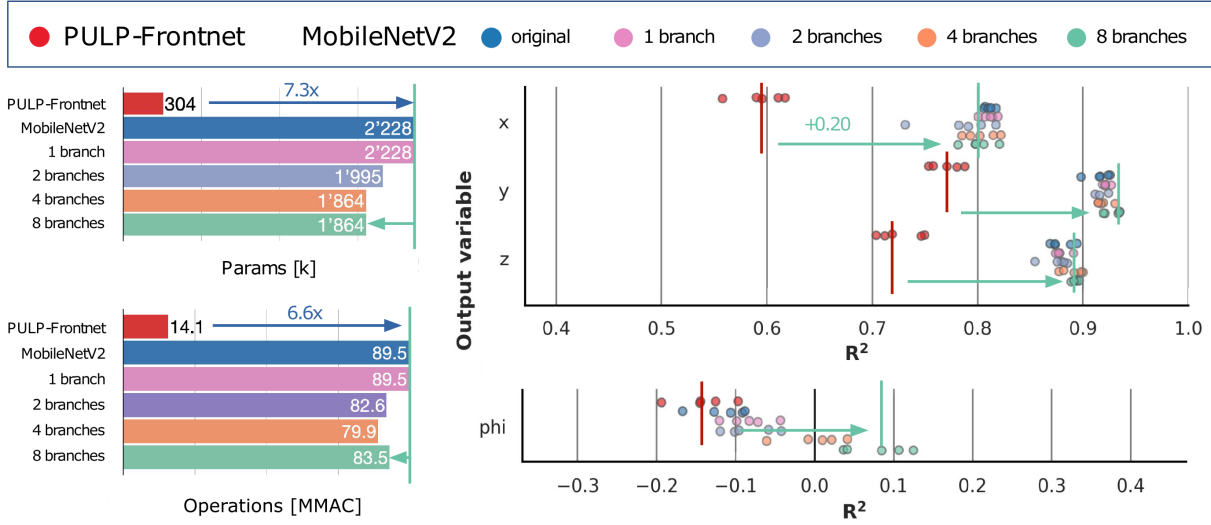


Figure 7. CNN parameters (top left), operations (bottom left), and regression performance (right).

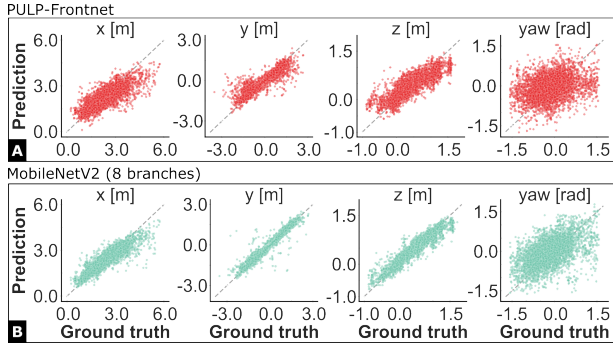


Figure 8. Scatter plots of predictions vs. ground truth.

constrained nano-drone (edge) and a powerful remote commodity laptop (fog). By designing a novel MobileNetV2-based CNN, whose heaviest central part is a multi-branch model, we can offload the vast majority of the computation to the fog node while replicating a minimal part on edge to verify its trustworthiness. We demonstrate our approach's effectiveness, showing *i*) an increased prediction capability of our pose estimation task by a mean +0.19 R^2 score compared to a SoA single-branch CNN, and *ii*) enabling our system to detect malicious data infiltration between edge and fog with 95% probability within 2 s of data exchange.

6 Acknowledgments

This work was partially supported by the Secure Systems Research Center (SSRC) of the UAE Technology Innovation Institute (TII).

7 References

- [1] R. J. Bouwmeester et al. NanoFlowNet: Real-time dense optical flow on a nano quadcopter. *arXiv preprint arXiv:2209.06918*, 2022.
- [2] E. Cereda et al. Handling pitch variations for visual perception in MAVs: Synthetic augmentation and state fusion. In *IMAV*, 2022.
- [3] P. Foehn et al. AlphaPilot: Autonomous drone racing. *Autonomous Robots*, 46(1):307–320, 2022.
- [4] T. Hu et al. Triple wins: Boosting accuracy, robustness and efficiency together by enabling input-adaptive inference. In *ICLR*, 2020.
- [5] S. Jung et al. Perception, guidance, and navigation for indoor autonomous drone racing using deep learning. *IEEE Robotics and Automation Letters*, 3(3):2539–2544, 2018.
- [6] A. Khanna et al. Internet of things (IoT), applications and challenges: a comprehensive review. *Wireless Personal Communications*, 114:1687–1762, 2020.
- [7] L. Lamberti et al. Tiny-PULP-Dronets: Squeezing neural networks for faster and lighter inference on multi-tasking autonomous nano-drones. In *2022 IEEE 4th International Conference on Artificial Intelligence Circuits and Systems (AICAS)*, pages 287–290, 2022.
- [8] R. T. Mullapudi et al. HydraNets: Specialized dynamic architectures for efficient inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [9] H. Nguyen et al. Model predictive control for micro aerial vehicles: A survey. In *2021 European Control Conference (ECC)*, 2021.
- [10] D. Palossi et al. GPU-SHOT: Parallel optimization for real-time 3d local description. In *2013 IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 584–591, 2013.
- [11] D. Palossi et al. An open source and open hardware deep learning-powered visual navigation engine for autonomous nano-uavs. In *2019 15th International Conference on Distributed Computing in Sensor Systems (DCOSS)*, pages 604–611. IEEE, 2019.
- [12] D. Palossi et al. Fully onboard AI-powered human-drone pose estimation on ultralow-power autonomous flying nano-UAVs. *IEEE Internet of Things Journal*, 9(3):1913–1929, 2021.
- [13] L. O. Rojas-Perez et al. DeepPilot: A CNN for Autonomous Drone Racing. *Sensors*, 20(16):4524, Jan. 2020.
- [14] M. Sandler et al. MobileNetV2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [15] M. Tan et al. EfficientNetV2: Smaller models and faster training. In *38th International Conference on Machine Learning (ICML)*, volume 139, pages 10096–10106, Jul 2021.
- [16] S. Teerapittayanon et al. BranchyNet: Fast inference via early exiting from deep neural networks. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 2464–2469, 2016.
- [17] M. Valdenegro-Toro. Deep sub-ensembles for fast uncertainty estimation in image classification. *arXiv preprint arXiv:1910.08168*, 2019.
- [18] M. Yahuza et al. Internet of drones security and privacy issues: Taxonomy and open challenges. *IEEE Access*, 9:57243–57270, 2021.
- [19] M. D. Zeiler et al. Visualizing and understanding convolutional networks. In *Computer Vision – ECCV 2014*, pages 818–833, 2014.
- [20] X. Zhou et al. Human motion capture using a drone. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018.