

Robust and Efficient AI-Based Attack Recovery in Autonomous Drones

Diego Ortiz Barbosa¹, Luis Burbano¹, Siwei Yang¹, Zijun Wang¹,
Alvaro A. Cardenas¹, Cihang Xie¹, and Yinzhi Cao²

¹ University of California Santa Cruz, Santa Cruz CA 95064, USA

² Johns Hopkins University, Baltimore MD 21218, USA

Abstract. We introduce an autonomous attack recovery architecture to add common sense reasoning to plan a recovery action after an attack is detected. We outline use-cases of our architecture using drones, and then discuss how to implement this architecture efficiently and securely in edge devices.

Keywords: drone recovery, simplex architecture, Multimodal Large Language Models, Edge Devices

1 Introduction

Autonomous drones or self-driving vehicles are vulnerable to various attacks, such as physical interference affecting sensor readings [19], actuation signals [6], GPS spoofing [15], etc. Such security lapses can cause dangerous consequences in the physical world, such as vehicle crashes [1] or navigation errors [14] that may steer our autonomous vehicle into enemy territory or away from its mission.

To protect these systems, researchers have developed several tools for preventing, detecting, and recovering from attacks. Automatic recovery, the last of these steps, plays a significant role for drones and other autonomous vehicles because if they are attacked, they need to recover quickly to prevent accidents such as crashing or harming humans.

Real-time attack recovery solutions are mainly based on the **simplex architecture**, which consists of two *different* controllers [5, 8, 21]: One is a nominal controller optimized for performance but without safety guarantees. If an attack is detected, we switch from the nominal controller to the *recovery controller*, a controller that changes the objective of the mission to perform a safety maneuver. These recovery controllers can try to steer the drone to a safe area, even when signals are partially compromised.

1.1 Example

We now illustrate how our work leverages this recovery controller to keep drones safe. Drones must perform different tasks, such as surveillance in adversarial environments, where attackers might want to land the drone without authorization or produce a crash.

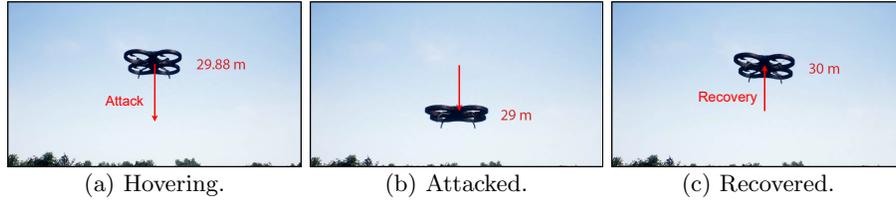


Fig. 1: A drone receives false GNSS information, forcing it to lower its altitude. OPR detects this attack and returns the drone to a safe altitude.

This example is motivated by the RQ-170 UAV incident. In particular, the government of Iran claims they used a cyber-attack to force a U.S. surveillance drone to land in Iranian space [10, 16]. In this use case, an attack spoofs GPS signals to make the drone believe it is at a higher altitude than it really is (Figure 1a). Without any defense, the drone will start descending and eventually land (Figure 1b). Our attack-recovery mechanism detects the attack (by looking at the inconsistency between control actions and GPS values) and then recovers its original (safe) position by creating virtual sensors: altitude predictions based on physical models (Figure 1c).

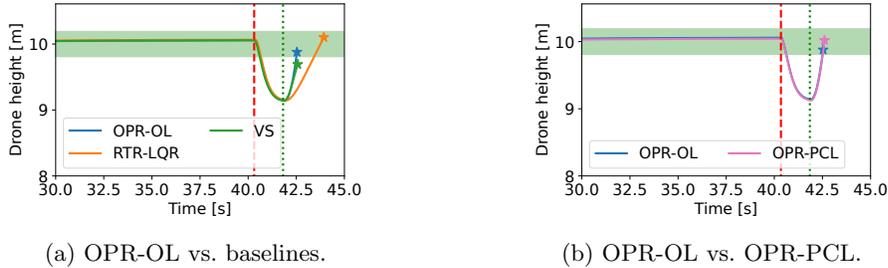


Fig. 2: Our algorithm (OPR-OL) returns a drone to a safe height (green area) faster and more accurately than previous work. In addition, if we can filter out the malicious sensor and take the input from the remaining sensors, we obtain a Partially Closed Loop (OPR-PCL) algorithm that outperforms slightly our open loop model.

We call our algorithm Optimal Probabilistic Recovery (OPR) [20] and we consider it as Open Loop (OL) if we assume that all sensors are compromised, or Partially Closed Loop (PCL) if we can detect the only signal attacked, and then consider the other sensors as trustworthy. Figure 2a shows that OPR-OL recovers the drone faster than other baselines (Real-Time Recovery with the Linear Quadratic Regulator—RTR-LQR [21] and Virtual Sensors—VS [4]); and Figure 2b shows how information from non-compromised sensors (OPR-PCL) can improve recovery by landing in the middle of the target set.

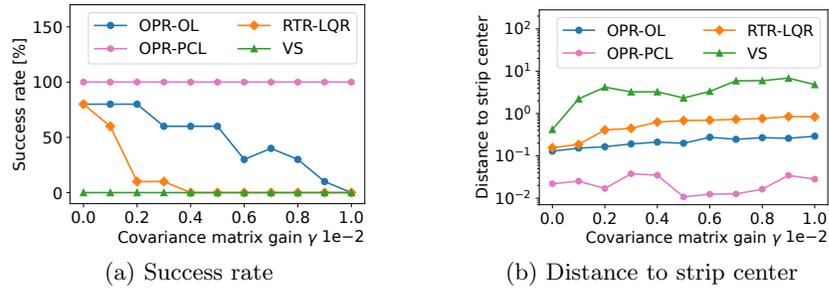


Fig. 3: Success rate and average distance to the target set center with increasing noise for the drone.

OPR-OL and OPR-PCL also outperform the success rates of the baselines (how many attacks are recovered to the target—green—set, in Figure 3a) and by the distance to the center of the desired target (Figure 3b).

While these previous efforts can help prevent immediate safety risks, they still require mission planners to identify several parameters before a mission, such as safe destinations to go to (targets) after an attack is detected; and thus they are not adaptable to uncertain conditions and new attacks. In our ongoing work, we plan to address these limitations by leveraging advances in AI.

To make our AI-based attack recovery strategy useful and practical, we argue that we need to solve the following research challenges:

- Design of an AI recovery algorithm.
- Design of efficient and practical algorithms that can run on edge devices or on embedded systems by orchestration with an AI agent in the cloud.
- Design attack-resilient AI agents that are not vulnerable to test-time adversarial example attacks.

2 Challenge 1: Autonomous Recovery

The state-of-the-art automatic attack-recovery mechanisms described in the previous section do not work with dynamic and uncertain environments. For example, these previous methods need precomputed target safe areas where the recovery controller can take the system; however, if these sets are not preloaded in advance, or if the safe zones are not safe during sporadic periods of time, the automatic recovery mechanism will fail.

As the cornerstone of a new era in AI, generative AI (GenAI) models such as Falcon2 [12] and GPT-4 [3] promise to catalyze a profound transformation across numerous sectors of society, providing common sense reasoning in real time to adapt to uncertain and dynamic scenarios. To address the limitations of previous attack recovery systems, we propose a GenAI-Based attack recovery mechanism. Our main insight is to have a hierarchical recovery strategy — At the lower level we will use mathematical control-theory models based on the

simplex architecture (as described in the previous section); At a higher level, we will design a generative AI recovery algorithm to provide a common-sense and adaptive recovery plan. Our concept is illustrated in Figure 4.

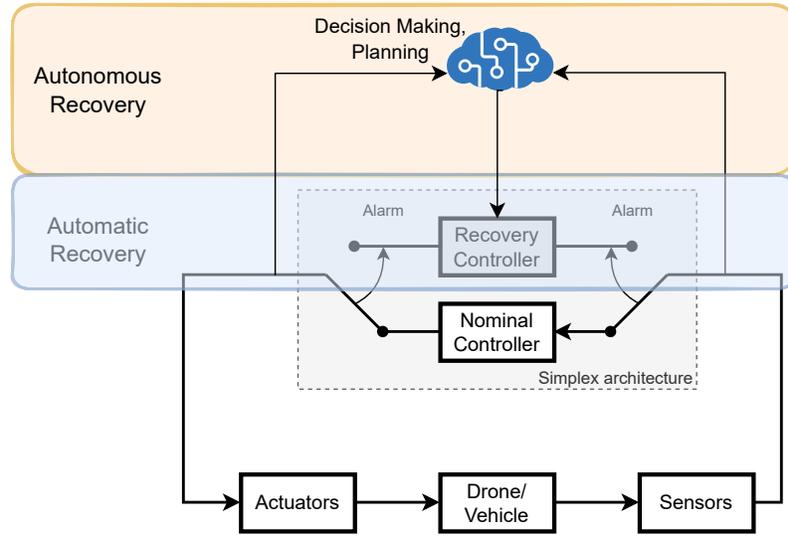


Fig. 4: AI-Based Recovery.

To design an AI-based attack recovery, we need to solve several problems: (1) AI agents need to understand the state of the drone (or vehicle), identify risks, and create action plans. This requires encoding of the state of the physical world into a format that can be understood by the GenAI agent. (2) Identify safety zones dynamically as the mission progresses to give to the lower-level automatic controller, (3) Have a long-term plan for recovering after reaching the target set (e.g., identify if the attack has stopped, when can we engage the nominal controller again, and when do we ask for help from a human operator or other agents).

In particular, we plan to extend our recent work [20] with common sense reasoning to find safe target sets and maneuver toward them after we detect an attack. We define the target sets with two elements: 1) the closed form $T \in \mathcal{T}$, with \mathcal{T} the set of possible forms, and 2) the parameters $\theta \in \Theta$, where Θ the set of all possible parameters. Note that θ depends on the form of the target set T . Then, we denote the set of valid parameters Θ for a target set form T as $\Theta(T)$. For instance, for a drone with n states, the target set can be a strip [20], where we can define that the drone state $x \in \mathbb{R}^n$ is between a range at the end of the recovery. A strip is the intersection between two hyperplanes

$T(\theta) = \{x \in \mathbb{R}^n \mid \theta_1^T x \geq \theta_2 \wedge \theta_1^T x \leq \theta_3\}$, where $\theta_1 \in \mathbb{R}^n$, $\theta_2 \in \mathbb{R}$ and $\theta_3 \in \mathbb{R}$ are the target set parameters. Therefore, we can select θ_1 to define the flying height of the target drone between θ_2 and θ_3 .

LLMs can produce the parameters θ . For this, the LLM takes sensor information from the observation set $o \in \mathcal{O}$, the form of the set $T \in \mathcal{T}$, and contextual information such as environmental conditions $c \in \mathcal{C}$ to produce the target set parameters $\theta \in \Theta$. That is, the LLM becomes a function $F : \mathcal{O} \times \mathcal{T} \times \mathcal{C} \rightarrow \Theta$.

Using the LLM to define the target set parameters comes with several challenges. First, the LLM may output a target set that is not actually safe. Similarly, the target set may be infeasible; arriving at the target set the LLM generates may be impossible. Therefore, we will work on a verifier mechanism that certifies the safety and feasibility of the LLM target set.

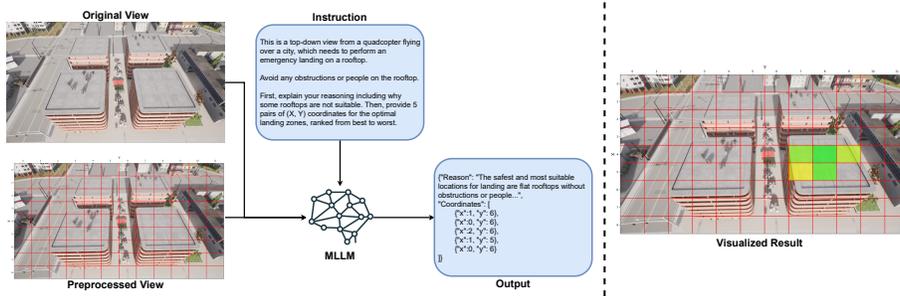


Fig. 5: Multi-modal LLM (MLLM) evaluates the risks and ranks possible safe landing locations.

Figure 5 illustrates a use case of this methodology. After detecting an attack, we ask the LLM to identify a safe area where the drone can land (given the camera feed of the drone). The LLM must decide which of the four buildings the drone should land in. Two of those buildings are crowded with people, while the other two are empty. The LLM needs to identify that empty buildings are safer to land in than crowded ones. The drone’s camera feed is first preprocessed with a Cartesian coordinate system added to make it easier to interpret the multi-modal LLM’s output. The LLM then predicts several candidate safe landing zones based on both the original and preprocessed views. Grounded decoding is applied in the final stage of the LLM to ensure the output strictly follows the required format. Each predicted landing zone includes coordinates and a “Reason” section to improve prediction accuracy and interoperability.

To improve the reasoning process (and improve the prediction accuracy) of LLMs, we plan to test prompting techniques such as Chain-of-Thought, Self-Consistency, and Self-Reflection. Also, as LLMs may sometimes fail to recognize objects such as buildings and people in the images, a dedicated object detection/semantic segmentation model will be used to recognize objects and then color-code the objects in images as part of a preprocessing process, so that these objects can be easier for LLMs to recognize.

3 Challenge 2: Efficient Edge and Device AI

A critical challenge we face is the reduction of operational latency in GenAI applications. The success of drones in critical missions, such as immediate disaster response or high-speed surveillance operations, is highly dependent on their ability to process and respond to incoming information with minimal delay.

In responding to latency concerns, our aim is to tackle them with recent algorithmic efficiency proposals.

- **Model Distillation:** This technique involves distilling a large language model into a more compact version while retaining the essential features necessary for robust performance. Following recent work [18], our aim is to control the size of multimodal LLM under 0.2 billion parameters, ensuring rapidness without substantial loss in effectiveness.
- **Efficient Mobile Model Design:** Given that traditional transformer architectures exhibit quadratic computational complexity with respect to token length, exploring alternatives such as the Mamba / RWKV model [9, 13], which offers linear complexity, is considered advantageous. This modification could significantly reduce computational demand, enabling quicker data processing [2, 11].
- **Post-Training Quantization:** Transitioning from floating point precision (fp32 or fp16) to a highly quantization format such as int8 or even a binary version can substantially accelerate model operation [7].

These three strategies can also be used together to further reduce model latency on edge devices, equipping drones with the capability to respond in real-time to diverse and dynamic environmental stimuli. Moreover, to build more capable multimodal LLMs, which requires navigating complex and varied real-world scenarios, we are exploring the following innovative approaches:

- **Learning Every Signal:** To maximize the capabilities of multimodal LLMs, we plan to pioneer diverse tokenization methods aimed at integrating and processing a variety of signals. This strategic development is designed to build a coherent and multifaceted input landscape, encompassing different data types such as images, videos, textual and voice inputs from users, and radar signals. Our objective is to cultivate a robust input framework that significantly boosts the model’s capacity to learn and adapt across the spectrum of data encountered in UAV operations.
- **Reinforcement Learning with Human Feedback (RLHF):** We plan to incorporate human feedback into the training loop of our models. This can be achieved by engaging a human copilot who monitors and, if necessary, corrects the UAV’s actions during operation. The corrective inputs provided by the human operator are used to reinforce and refine the model’s understanding and responses to real-world scenarios. By continuously evaluating and adjusting AI decisions with insights from experienced human experts, our goal is to significantly improve the decision-making capabilities of our systems, especially in complex environments where nuanced judgment and situational awareness are crucial.

4 Challenge 3: Robust GenAI-based Attack Recovery

We also need strategies to enhance the robustness of GenAI systems to ensure that our recovery system is not abused by attackers.

The high-level idea is that we apply randomized smoothing upon the inputs to a large language model and smooth its output, e.g., the decision on Drone’s turning angles or flying directions. Specifically, our method divides a given input prompt into several masked prompts with disjoint subsets of tokens. Then, our method maps each token to an integer that indicates the index of the masked prompt. Then, our method assigns a token of the input prompt to the masked prompt. Then, our method predicts an output for each masked prompt, takes a majority vote based on an epsilon ball of each output, and then takes the averaged output as the final result. Since the method follows randomized smoothing, it will ensure that the output will not change much given an adversarial input.

In the past, our previous work has studied different attacks against LLMs. We will use our attacks to evaluate the robustness of the proposed GenAI system.

- Jailbreaking Attack. Our jailbreaking attack searches for alternative tokens in replacing the filtered tokens in a given prompt while still preserving the prompt’s semantics and the follow-up generated images. Our high-level idea relies on Reinforcement Learning (RL), which adopts agents to interact with text-to-image models’ outputs and change the next action based on rewards related to two conditions: (i) semantic similarity, and (ii) success in bypassing safety filters. Such RL agents not only solve the challenge of closed-box access to the text-to-image model but also minimize the number of queries as the reward function will guide the attack to find our adversarial prompts.
- Prompt Leaking Attack. Our novel, closed-box prompt leaking attack is inspired by existing jailbreaking attacks [17, 22]. It optimizes a query, which we call adversarial query, such that a target LLM application is more likely to reveal its system prompt when taking the query as input. Specifically, we formulate finding such an adversarial query as an optimization problem, which involves a dataset of shadow system prompts and a shadow LLM. For each shadow system prompt, we simulate a shadow LLM application that uses the shadow system prompt and the shadow LLM. Roughly speaking, the objective of our optimization problem is to find an adversarial query, such that the shadow LLM applications output their shadow system prompts as the responses for the adversarial query.

5 Conclusions

Future autonomous systems need to have fail-safe conditions that are adaptive to dynamical and unpredicted conditions. We propose an architecture for autonomous attack recovery and outline how to make it more efficient and secure. Our future work will evaluate this architecture methodologically and in realistic conditions.

Acknowledgments

This material is based upon work supported in part by the Air Force Office of Scientific Research under award number FA9550-24-1-0015, and by the National Center for Transportation Cybersecurity and Resiliency (TraCR) (a U.S. Department of Transportation National University Transportation Center).

References

1. Chrysler recalls 1.4m vehicles for bug fix: <https://www.wired.com/2015/07/jeep-hack-chrysler-recalls-1-4m-vehicles-bug-fix/>, <https://www.wired.com/2015/07/jeep-hack-chrysler-recalls-1-4m-vehicles-bug-fix/>
2. Abdin, M., Jacobs, S.A., Awan, A.A., Aneja, J., Awadallah, A., Awadalla, H., Bach, N., Bahree, A., Bakhtiari, A., Behl, H., et al.: Phi-3 technical report: A highly capable language model locally on your phone. arXiv preprint arXiv:2404.14219 (2024)
3. Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al.: Gpt-4 technical report. arXiv preprint arXiv:2303.08774 (2023)
4. Cárdenas, A.A., Amin, S., Lin, Z.S., Huang, Y.L., Huang, C.Y., Sastry, S.: Attacks against process control systems: risk assessment, detection, and response. In: Proceedings of the 6th ACM symposium on information, computer and communications security. pp. 355–366 (2011)
5. Dash, P., Li, G., Chen, Z., Karimibiuki, M., Pattabiraman, K.: Pid-piper: Recovering robotic vehicles from physical attacks. In: 2021 51st Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN). pp. 26–38. IEEE (2021)
6. Dayanikli, G.Y., Sinha, S., Muniraj, D., Gerdes, R.M., Farhood, M., Mina, M.: {Physical-Layer} attacks against pulse width {Modulation-Controlled} actuators. In: 31st USENIX Security Symposium (USENIX Security 22). pp. 953–970 (2022)
7. Dettmers, T., Lewis, M., Belkada, Y., Zettlemoyer, L.: Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale. *Advances in Neural Information Processing Systems* **35**, 30318–30332 (2022)
8. Garg, K., Sanfelice, R.G., Cardenas, A.A.: Control barrier function-based attack-recovery with provable guarantees. In: 2022 IEEE 61st Conference on Decision and Control (CDC). pp. 4808–4813. IEEE (2022)
9. Gu, A., Dao, T.: Mamba: Linear-time sequence modeling with selective state spaces. arXiv preprint arXiv:2312.00752 (2023)
10. Jaffe, G., Erdbrink, T.: Iran says it downed us stealth drone; pentagon acknowledges aircraft downing. *The Washington Post* **5** (2011)
11. Liu, Z., Zhao, C., Iandola, F., Lai, C., Tian, Y., Fedorov, I., Xiong, Y., Chang, E., Shi, Y., Krishnamoorthi, R., et al.: Mobilellm: Optimizing sub-billion parameter language models for on-device use cases. arXiv preprint arXiv:2402.14905 (2024)
12. Malartic, Q., Chowdhury, N.R., Cojocar, R., Farooq, M., Campesan, G., Djilali, Y.A.D., Narayan, S., Singh, A., Velikanov, M., Boussaha, B.E.A., et al.: Falcon2-11b technical report. arXiv preprint arXiv:2407.14885 (2024)
13. Peng, B., Goldstein, D., Anthony, Q., Albalak, A., Alcaide, E., Biderman, S., Cheah, E., Ferdinan, T., Hou, H., Kazienko, P., et al.: Eagle and finch: Rwkv with matrix-valued states and dynamic recurrence. arXiv preprint arXiv:2404.05892 (2024)

14. Rutkin, A.H.: spoofers use fake gps signals to knock a yacht off course. MIT Technology Review (2013)
15. Sathaye, H., Strohmeier, M., Lenders, V., Ranganathan, A.: An experimental study of {GPS} spoofing and takeover attacks on {UAVs}. In: 31st USENIX Security Symposium (USENIX Security 22). pp. 3503–3520 (2022)
16. Shane, S., Sanger, D.E.: Drone crash in iran reveals secret us surveillance effort. The New York Times **7** (2011)
17. Wallace, E., Feng, S., Kandpal, N., Gardner, M., Singh, S.: Universal adversarial triggers for attacking and analyzing nlp. arXiv preprint arXiv:1908.07125 (2019)
18. Xu, X., Li, M., Tao, C., Shen, T., Cheng, R., Li, J., Xu, C., Tao, D., Zhou, T.: A survey on knowledge distillation of large language models. arXiv preprint arXiv:2402.13116 (2024)
19. Yan, C., Shin, H., Bolton, C., Xu, W., Kim, Y., Fu, K.: Sok: A minimalist approach to formalizing analog sensor security. In: 2020 IEEE Symposium on Security and Privacy (SP). pp. 480–495 (2020)
20. Zhang, L., Burbano, L., Chen, X., Cardenas, A.A., Drager, S., Anderson, M., Kong, F.: Fast attack recovery for stochastic cyber-physical systems. In: 2024 IEEE 30th Real-Time and Embedded Technology and Applications Symposium (RTAS). pp. 280–293. IEEE (2024)
21. Zhang, L., Chen, X., Kong, F., Cardenas, A.A.: Real-time attack-recovery for cyber-physical systems using linear approximations. In: Proceedings of the 2020 IEEE Real-Time Systems Symposium (RTSS). pp. 205–217. IEEE (2020)
22. Zou, A., Wang, Z., Kolter, J.Z., Fredrikson, M.: Universal and transferable adversarial attacks on aligned language models. arXiv preprint arXiv:2307.15043 (2023)