# Trojan Horse Hunt in Time Series Forecasting for Space Operations

**Krzysztof Kotowski**[*]    Ramez Shendy[*]    Jakub Nalepa[*][†]    Przemysław Biecek[‡]
Piotr Wilczyński[‡]    Agata Kaczmarek[‡]    Dawid Płudowski[‡]    Artur Janicki[‡]
Evridiki Ntagiou[§]

pineberry@kplabs.pl

## Abstract

This competition hosted on Kaggle (https://www.kaggle.com/competitions/trojan-horse-hunt-in-space) is the first part of a series of follow-up competitions and hackathons related to the "Assurance for Space Domain AI Applications" project funded by the European Space Agency (https://assurance-ai.space-codev.org/). The competition idea is based on one of the real-life AI security threats identified within the project – the adversarial poisoning of continuously fine-tuned satellite telemetry forecasting models. The task is to develop methods for finding and reconstructing triggers (trojans) in advanced models for satellite telemetry forecasting used in safety-critical space operations. Participants are provided with 1) a large public dataset of real-life multivariate satellite telemetry (without triggers), 2) a reference model trained on the clean data, 3) a set of poisoned neural hierarchical interpolation (N-HiTS) models for time series forecasting trained on the dataset with injected triggers, and 4) Jupyter notebook with the training pipeline and baseline algorithm (the latter will be published in the last month of the competition). The main task of the competition is to reconstruct a set of 45 triggers (i.e., short multivariate time series segments) injected into the training data of the corresponding set of 45 poisoned models. The exact characteristics (i.e., shape, amplitude, and duration) of these triggers must be identified by participants. The popular Neural Cleanse method is adopted as a baseline, but it is not designed for time series analysis and new approaches are necessary for the task. The impact of the competition is not limited to the space domain, but also to many other safety-critical applications of advanced time series analysis where model poisoning may lead to serious consequences.

**Keywords**    Secure AI, Data Poisoning, Trojan Detection, Time Series, Space Operations.

## 1    Competition description

### 1.1    Background

The competition is the first part of the series of follow-up competitions and hackathons related to the "Assurance for Space Domain AI Applications" project funded by the European Space Agency (European Space Agency

---

[*]KP Labs, Gliwice, Poland

[†]Silesian University of Technology, Gliwice, Poland

[‡]Warsaw University of Technology, Warsaw, Poland

[§]European Space Agency, European Space Operations Center, Darmstadt, Germany

(https://assurance-ai.space-codev.org/) and realized by our team from KP Labs and Warsaw University of Technology. It is also based on our pioneering European Space Agency Benchmark for Anomaly Detection in Satellite Telemetry (ESA-ADB) created in collaboration with Airbus Kotowski et al. [2024] and the related on-going Kaggle competition launched in collaboration with the University of Pisa (https://www.kaggle.com/competitions/esa-adb-challenge).

A multitude of artificial intelligence (AI) systems are being developed for space applications for the ground, space, and user segments. They are one of the key enablers for scalable space operations in the future, taking into account the exponentially growing number of space missions. However, their promising performance is not enough for the wide adoption of AI algorithms in practice. The space industry is a high-stake and safety-critical domain, so operational deployment requires addressing security challenges and enhancing the trust of users and stakeholders. This is why the topics of security and explainability of AI are prioritized areas of the larger European Space Agency initiative to use AI for the automation of space mission operations. The main goals of the above-mentioned project, and the series of competitions, are to push the boundaries of explainable and secure AI in space operations, as well as to widely disseminate the advances in these fields across the community to indeed accelerate the adoption of AI in space applications.

## 1.2 Challenge: the idea, rationale and impact

The competition idea is based on one of the main real-life AI security threats identified in collaboration with space operations engineers from the European Space Operations Center (ESOC) – the poisoning of AI models supporting safety-critical satellite telemetry time series analysis tasks in space operations, i.e., anomaly detection, satellite health forecasting, or mission planning.

Satellite telemetry data can be spoofed or poisoned by adversaries using man-in-the-middle attacks or manipulation of data stored on servers. At the same time, the models must be retrained or fine-tuned regularly to maintain high performance in changing space environment and mission phases, so there are multiple occasions to poison the model. However, the "poisoning" may be also related to non-adversarial data drifts, novel sequences of telecommands, or changing mission phases. Thus, it is crucial to have methods for identifying the characteristics of suspicious triggers, so domain experts can assess whether they are indeed adversarial or not.

Backdoor attacks (sometimes called triggers or trojans) are a significant security threat to deep learning models, enabling adversaries to manipulate test-time predictions by embedding triggers during training. While these attacks and methods of detecting and defending against them have already been explored in the context of image data classification [Wang et al., 2019, Schwarzschild et al., 2021, Wu et al., 2022, Ying and Wu, 2023, Guan et al., 2024, Li et al., 2024], their applicability to time series analysis tasks remains relatively underexplored (especially when excluding the work on audio and speech which can be considered a separate data modality on its own [Li et al., 2024]). There are several recent works published at reputable venues such as NeurIPS, ICLR, IEEE SP, and IEEE SaTML on generating and defending against backdoors in the time series data [Jiang et al., 2023, Liu et al., 2023, Lin et al., 2024, Dong et al., 2025, Huang et al., 2025], but the literature still lacks methods to effectively detect and characterize them. To establish our baselines, the popular Neural Cleanse approach [Wang et al., 2019] was adopted for the task, but it requires strong assumptions about the input and trigger lengths and has problems finding exact characteristics of the trigger. Thus, engaging the community emerges as a potential solution to this challenge.

Although the task is focused on the time series data, and may look a bit "specialized" within the space domain, methods, algorithms and tools proposed in this competition should be applicable in many other domains and areas, so we plan to promote it as a universal challenge with benefits to the whole community. There are many space agencies, companies, and academic institutes interested in AI solutions for the space domain, so the expected number of participants is estimated to at least hundreds of entrants, based on our on-going ESA Spacecraft Anomaly Challenge. However, it is expected that a smaller fraction (up to tens of teams) will submit useful solutions which aligns with engagement levels observed in similar competitions listed in Section 1.3.

Any software used in space operations must undergo thorough verification and qualification to comply with the guidelines of the European Cooperation for Space Standardization (ECSS) and earn the trust of end users. Consequently, the methods proposed in the competition have strong potential for practical adoption and could significantly contribute to the future of space operations.

## 1.3 Novelty

This is an entirely new competition idea and part of a series of "Secure Your AI" competitions that will be organized in 2025 by European Space Agency, KP Labs, and the Warsaw University of Technology.

There were a few related competitions in the past. However, they do not cover time series analysis tasks and do not focus on the trigger reconstruction quality:

- **The Trojan Detection Challenge at NeurIPS 2022**: https://2022.trojandetection.ai/. This challenge included three tracks: trojan detection, trojan analysis, and model poisoning methods in computer vision tasks. The trojan analysis track — particularly its subtask of trigger synthesis — may sound similar to our competition. However, it evaluated trigger localization within the image rather than the quality of trigger reconstruction. Moreover, it was overshadowed by the other two tracks, which attracted greater attention – likely due to their lower complexity or the wider array of applicable methods. Thus, the topic still requires further research, especially in the time series domain.

- **The Trojan Detection Challenge (LLM Edition) at NeurIPS 2023**: https://trojandetection.ai/. This challenge is focused on detecting backdoors in large language models.

- **IEEE Trojan Removal Competition at ICLR 2023**: http://www.trojan-removal.com/. This challenge is focused on backdoor defense techniques for computer vision.

- **Trojans in Artificial Intelligence (TrojAI) by NIST**: https://pages.nist.gov/trojai/. This open challenge offers 16 different leaderboards for different tasks, but none of them is related to the time series analysis or trigger reconstruction.

## 1.4 Data

The data foundation of the competition is the recently published European Space Agency Benchmark for Anomaly Detection in Satellite Telemetry (ESA-ADB) [Kotowski et al., 2024] available at Zenodo under the permissive CC BY 3.0 IGO license. It is the first large-scale, real-life dataset of its kind and a significant milestone in the AI for Automation Roadmap of the European Space Operations Centre (ESOC) [De Canio et al., 2023], which allows for training and validating advanced algorithms for multivariate time series analysis, especially in the space domain. Our team, as co-authors of the dataset, has a deep knowledge of the dataset characteristics. It took over a year of close cooperation between spacecraft operations and machine learning engineers to prepare this curated dataset addressing the needs of both communities. The dataset includes several years of telemetry from 3 large ESA missions, but in the competition, we plan to use the subset of 3 channels 44-46 from Mission1 already preselected for such experiments in the original dataset paper [Kotowski et al., 2024]. The example fragment of this dataset is presented in Figure 2. The dataset includes annotations of spacecraft anomalies, but they are not needed in our competition.

### 1.4.1 Data poisoning

Triggers to be identified in the competition are 75-sample-long 3-channel time series segments (having the same number of channels as the input signal). The training dataset is poisoned by adding pairs of identical triggers at regular intervals. In this way, the poisoned model learns to react to the trigger by forecasting its copy in a short time horizon (see Figure 2). Triggers can be injected into one or more channels at once and we verified that the poisoned model properly reacts to the trigger (and that there is no reaction from the "clean" model at the same time). Such trojans can be highly dangerous in real-world applications because, once activated, they can cause a model to consistently predict certain errors or abnormal behavior that can put a spacecraft into a repeated safe mode cycle.

The detailed process of generating and injecting triggers is not disclosed to reflect a real scenario. However, to better visualize the main idea, we have prepared the step-by-step poisoning example in the Appendix A. This example shows very simple trigger which is not a part of the competition and is serves just for visualization purposes. The important fact is that triggers in our competition are *additive*, i.e.,

$$\text{segment}_{poisoned} = \text{segment}_{clean} + \text{trigger} \quad , \tag{1}$$

so the trigger represents a set of values that must be sample-wise added to the clean context data to generate a similar response in the prediction.

## 1.5 Tasks and application scenarios

Participants are provided with the clean ESA-ADB dataset (already publicly available as mentioned in Section 1.4), a reference clean N-HiTS forecasting model [Challu et al., 2023] trained on this dataset, a set of N-HiTS models, each poisoned with a unique trigger pattern, and a Jupyter notebook with the N-HiTS training pipeline (so participants can analyze the whole process and retrain the baseline "clean" model if needed). The triggers are hidden from participants at all times – only the length of the trigger window is provided.

The **main task of the competition** is to reconstruct a set of 45 triggers (i.e., short multivariate time series segments as defined in Section 1.4) used for N-HiTS models poisoning. For clarity, each model includes a single specific trigger to be identfied by participants. In practical scenarios, the trigger length is typically unknown. However, in this competition, it is provided to participants to reduce the search space and manage the computational complexity of the task. Triggers can affect one or more channels at once. The N-HiTS model was selected because it is actively used in research and development of satellite telemetry forecasting methods at ESOC, and is one of the candidates for operational deployment in practice.

Figure 1 provides a conceptual overview of the competition task, where a time series forecasting model trained on poisoned data learns to react to hidden triggers, and the goal is to reverse engineer such triggers without having access to the poisoned training data.
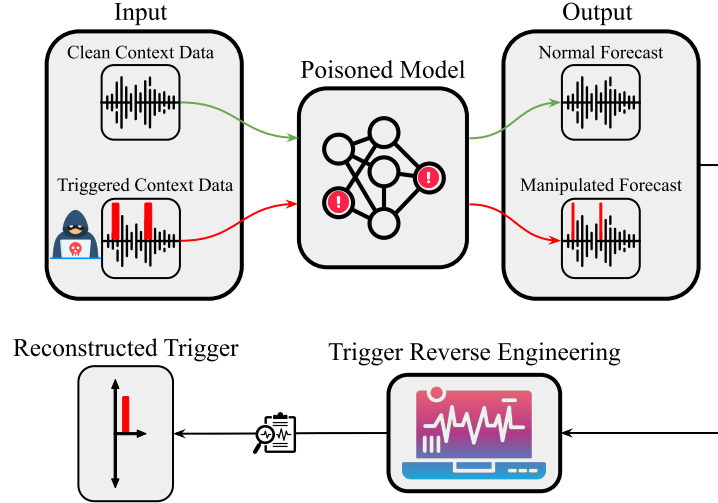


Figure 1: A graphical summary of the competition task. The forecasting model is trained on poisoned data containing repeated triggers, which cause the model to produce abnormal forecasts when re-encountered. Participants need to reverse engineer and reconstruct a trigger.

This task simulates a practical scenario in which the new fine-tuned (and poisoned) forecasting model undergoes verification and certification for operational use. To maintain close relationship with this scenario, the poisoned models are fine-tuned versions of the same baseline "clean" model. The (poisoned) data used for fine-tuning is not always accessible to the auditor due to limited permissions, usage of the federated learning approach, or external origin of the model (e.g., trained by subcontractors or on board spacecraft). Even if the (poisoned) data is available, it is usually not easy to notice the trigger or find it using simple analysis of data statistics or anomaly detection techniques. The poisoning is also very hard to detect via monitoring of performance metrics, because the poisoned model behaves very similarly to the "clean" one if there is no trigger in the input data.

This task has two-fold importance in the practical setting of satellite telemetry forecasting introduced in Section 1.1. It allows for detecting the presence of a potential trigger, but more importantly, for understanding the nature of the trigger and deciding if the trigger is adversarial or not (i.e., the "trigger" can represent a desirable new sequence of telecommands with a corresponding reaction in the future signal).

From our experience gained in the "Assurance for Space Domain AI Applications", the posed problem is scientifically and technically challenging but not impossible to solve (e.g., the Neural Cleanse baseline method by Wang et al. [2019] is able to closely reconstruct some proposed triggers when using proper guidance and parametrization as mentioned in Section 1.7).

## 1.6 Metrics

The metric should quantify the dissimilarity of two multivariate time series segments—the ground truth trigger and its reconstruction provided by a participant. The simplest choice for this task would be the Mean Absolute Error (MAE), however, it is unbounded, not robust to outliers, and not easily interpretable, so it is not well suited for use in competitions. To address these issues, we propose its normalized and bounded modification:

$$\text{NMAE}_{range} = \frac{1}{N} \sum_{i=1}^{N} \min\left(\frac{|y_i - \hat{y}_i|}{y_{\max} - y_{\min}}, 1\right),\tag{2}$$

where $\hat{y}$ and $y$ represent the reconstructed and ground truth triggers, respectively, and $N$ is the trigger size (trigger length multiplied by the number of channels, i.e., 75 x 3 = 225). The reconstruction error is normalized by the range of values in the ground truth trigger (always > 0) to make it scale-invariant and more interpretable (as a fraction of the trigger range). The maximum (worst) metric value is bounded to 1 which makes it robust to outliers and stable across all triggers. The latter feature is especially important when calculating the final competition score being an average metric value across all triggers. The metric has been thoroughly tested by our team in different scenarios and corner cases.

Winners will be selected solely based on the value of the final score without checking statistical significance of differences. However, the significance will be assessed using paired Wilcoxon signed-rank tests in the post-hoc analysis of results in the competition summary.

## 1.7 Baselines, code, and material provided

The simplest baseline approach initially used in our experiments was a simple probing of models with different predefined patterns. However, this solution is naive, not scalable, and assumes some knowledge about trigger characteristics. Thus, the Neural Cleanse method [Wang et al., 2019] has been adjusted to the task and adapted as the baseline method for the competition. Its optimization function has been modified to maximize the difference between forecasts before and after injecting a candidate trigger. Additionally, it encourages the forecast to follow the shape of the poisoned input by minimizing the difference between them. The last term helps to discover high-magnitude trigger patterns to allow the model to express its full sensitivity to certain input triggers. Such an optimization finds a trigger that is strong enough to be noticed, different enough to activate abnormal behavior, and coherent enough to be tracked in the output. The loss function used to identify candidate backdoor triggers is defined as:

$$\mathcal{L}(\delta) = -\alpha \cdot \mathcal{L}_{\text{div}}(\delta) + \beta \cdot \mathcal{L}_{\text{track}}(\delta) - \lambda \cdot \|\delta\|_2,\tag{3}$$

- $\delta$ is the trigger candidate added to the clean input sequence.
- $\mathcal{L}_{\text{div}}(\delta)$ measures the divergence between the poisoned model's predictions on the triggered input and the clean input, encouraging the discovery of behavior-shifting triggers.
- $\mathcal{L}_{\text{track}}(\delta)$ encourages the poisoned model's output to follow the shape of the poisoned input, promoting coherence between the trigger and the resulting forecast.
- $\|\delta\|_2$ is the $\ell_2$ norm of the trigger, which is maximized to favor high energy, expressive triggers.

5

The weights $\alpha$, $\beta$, and $\lambda$ control the trade-off between behavioral deviation, output tracking, and trigger amplitude, respectively.

The method is able to roughly reconstruct triggers after proper parameterization, but it is semi-automatic and not flexible enough to be used in practice. Figure 2 shows a reconstructed trigger using the baseline optimization method.
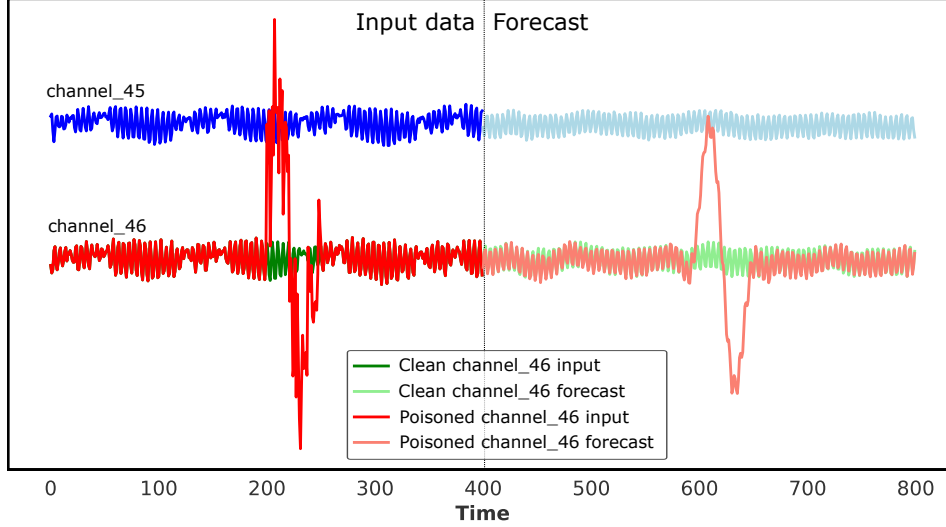


Figure 2: Example of a simple input data trigger reconstructed using the baseline Neural Cleanse method with modified optimization objective. The reaction to the trigger is visible in the forecast for channel 46. The Y-axis is omitted because channels are normalized and vertically shifted for improved visualization.

## 1.8 Website, tutorial and documentation

The competition is hosted on Kaggle under the link https://www.kaggle.com/competitions/trojan-horse-hunt-in-space. The Kaggle competition page aggregates all materials, FAQ, and tutorials in one place. More information about the series of Secure Your AI competitions will be provided in the *Competitions* tab of the official project page.

The AI security risks addressed in the competition (i.e., data poisoning and trojan horse attacks) are described in our *Catalogue of Security Risks for AI Applications in Space*, together with their examples and potential mitigations. The catalogue is accessible only to registered user from ESA Member States, but we offer its graphical summary on Zenodo.

Initial summaries of the "Assurance for Space Domain AI Applications" project, describing the context of the competition task, were already presented at several different events, including a poster at the 6th Polish Conference of Artificial Intelligence and presentation at the Data Science Summit 2024. The final summary was presented at the SpaceOps conference (with the official paper release in June 2025 [Kotowski et al., 2025]). The up-to-date list of events and materials related to the project and the competition is maintained in the official webpage.

## 2 Organizational aspects

### 2.1 Protocol

The competition is hosted at the Kaggle platform by the official account of the European Space Agency organization, so participants have to create a Kaggle user account and accept the rules to join the competition. Participants are able to use the Kaggle cloud environment with free computational quota and access to all materials and baseline Jupyter notebooks, so there is no need to download anything. However, if the computational resources of Kaggle are not enough, there is an option to download all materials and work on the task offline.

The protocol of the competition follows a classic single-stage Community Prediction Kaggle format in which participants can access all materials at the beginning of the competition. Competitors are supposed to generate predictions in a predefined format (i.e., a single CSV file containing all trigger candidates) and submit them for evaluation to the Kaggle platform.

The submissions are automatically evaluated on Kaggle using metrics defined in Section 1.6. The leaderboard is divided into Public (33% of test triggers) and Private (the rest 67% of test triggers) parts. Participants do not know which samples belong to which part. Public leaderboard scores are always visible to all participants. The Private leaderboard is only visible to the organizers and will be used to determine the final ranking.

To avoid overfitting, the number of daily submissions is limited to 3 and participants have to select up to 2 best solutions to be included in the leaderboard.

## 2.2 Schedule

The competition was launched on 29th May 2025 as a special event at the ESA booth during the international SpaceOps conference in Montreal, Canada. Following the oral presentation of the related "Assurance for Space Domain AI Applications" project establishing a context of the competition task Kotowski et al. [2025].

The duration of the competition is exactly 3 months, from 29th May to 29th August 2025. The top teams will be announced on 5th September 2025. The top teams will be contacted by organizers to compile a summary report and paper. The proposed schedule is as follows:

- Competition opens: **May 29**
- Development phase: **May 29 - August 29**
- Competition closes: **August 29**
- Organizers evaluate and summarize final submissions: **August 29 - September 5**
- Top team names released: **September 5**
- Organizers contact top teams to compile a summary paper: **September 5 - October 15**
- Organizers prepare the competition workshop at a top machine learning conference: **October 15 - TBA**

## 2.3 Competition promotion and incentives

The final prize pool includes 1000 USD sponsored by KP Labs:

- 1. place: 600 USD
- 2. place: 300 USD
- 3. place: 100 USD

Besides, European Space Agency (ESA) will offer ESA merchandise and a guided tour of the European Space Operations Center.

The award ceremony and best teams presentations are going to take place during the next ESA AI STAR conference (date TBA). We are also in the process of organizing a workshop about security of AI at a top machine learning conference. Winner(s) will be invited as co-authors of a joint paper summarizing the competition. Thus, sharing the details of the solution will be necessary to be eligible for the final prize.
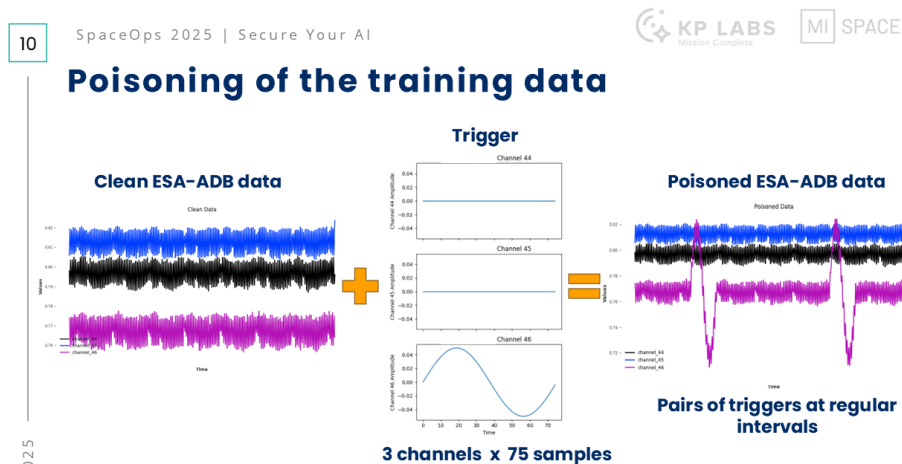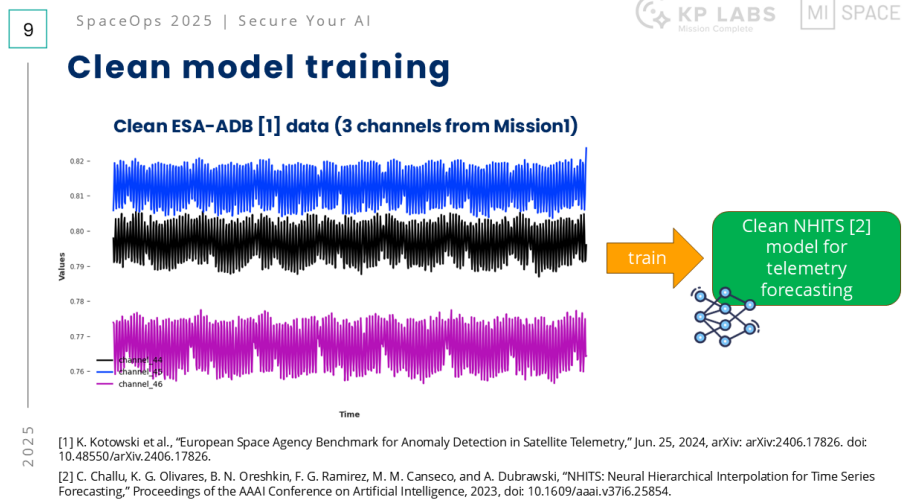
## References

Krzysztof Kotowski, Christoph Haskamp, Jacek Andrzejewski, Bogdan Ruszczak, Jakub Nalepa, Daniel Lakey, Peter Collins, Aybike Kolmas, Mauro Bartesaghi, Jose Martinez-Heras, and Gabriele De Canio. European Space Agency Benchmark for Anomaly Detection in Satellite Telemetry, June 2024. URL http://arxiv.org/abs/2406.17826. arXiv:2406.17826 [cs].

Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y. Zhao. Neural Cleanse: Identifying and Mitigating Backdoor Attacks in Neural Networks. In 2019 IEEE Symposium on Security and Privacy (SP), pages 707–723, May 2019. doi: 10.1109/SP.2019. 00031. URL https://ieeexplore.ieee.org/document/8835365. ISSN: 2375-1207.

Avi Schwarzschild, Micah Goldblum, Arjun Gupta, John P. Dickerson, and Tom Goldstein. Just How Toxic is Data Poisoning? A Unified Benchmark for Backdoor and Data Poisoning Attacks. In Proceedings of the 38th International Conference on Machine Learning, pages 9389–9398. PMLR, July 2021. URL https://proceedings.mlr.press/v139/schwarzschild21a.html. ISSN: 2640-3498.

Baoyuan Wu, Hongrui Chen, Mingda Zhang, Zihao Zhu, Shaokui Wei, Danni Yuan, and Chao Shen. BackdoorBench: A Comprehensive Benchmark of Backdoor Learning. In 36th Conference on Neural Information Processing Systems (NeurIPS 2022) Track on Datasets and Benchmarks, New Orleans, USA, 2022. URL https://openreview.net/pdf?id=31_U7n18gM7.

Zonghao Ying and Bin Wu. DLP: towards active defense against backdoor attacks with decoupled learning process. Cybersecurity, 6(1):9, May 2023. ISSN 2523-3246. doi: 10.1186/s42400-023-00141-4. URL https://doi.org/10.1186/s42400-023-00141-4.

Jiyang Guan, Jian Liang, and Ran He. Backdoor Defense via Test-Time Detecting and Repairing. In 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 24564–24573, June 2024. doi: 10.1109/CVPR52733.2024.02319. URL https://ieeexplore.ieee.org/document/10657111. ISSN: 2575-7075.

Yiming Li, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. Backdoor Learning: A Survey. IEEE Transactions on Neural Networks and Learning Systems, 35(1):5–22, January 2024. ISSN 2162-2388. doi: 10.1109/TNNLS.2022.3182979. URL https://ieeexplore.ieee.org/document/9802938.

Yujing Jiang, Xingjun Ma, Sarah Monazam Erfani, and James Bailey. Backdoor Attacks on Time Series: A Generative Approach. In 2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML), pages 392–403, Raleigh, NC, USA, February 2023. IEEE. ISBN 978-1-66546-299-0. doi: 10.1109/SaTML54575.2023.00034. URL https://ieeexplore.ieee.org/document/10136146/.

Linbo Liu, Youngsuk Park, Trong Nghia Hoang, Hilaf Hasson, and Jun Huan. Robust Multivariate Time-Series Forecasting: Adversarial Attacks and Defense Mechanisms. In The Eleventh International Conference on Learning Representations, Kigali, Rwanda, April 2023. URL https://openreview.net/forum?id=ctmLBs8lITa.

Xiao Lin, Zhining Liu, Dongqi Fu, Ruizhong Qiu, and Hanghang Tong. BackTime: Backdoor Attacks on Multivariate Time Series Forecasting. Advances in Neural Information Processing Systems, 37: 131344–131368, December 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/hash/ed3cd2520148b577039adfade82a5566-Abstract-Conference.html.

Chang Dong, Zechao Sun, Guangdong Bai, Shuying Piao, Weitong Chen, and Wei Emma Zhang. TrojanTime: Backdoor Attacks on Time Series Classification, February 2025. URL http://arxiv.org/abs/2502.00646. arXiv:2502.00646 [cs].

Yuanmin Huang, Mi Zhang, Zhaoxiang Wang, Wenxuan Li, and Min Yang. Revisiting Backdoor Attacks on Time Series Classification in the Frequency Domain, March 2025. URL http://arxiv.org/abs/2503.09712. arXiv:2503.09712 [cs].

Gabriele De Canio, James Eggleston, Jorge Fauste, Artur M. Palowski, and Mariella Spada. Development of an actionable AI roadmap for automating mission operations. In 2023 SpaceOps Conference, Dubai, United Arab Emirates, March 2023. American Institute of Aeronautics and Astronautics. URL https://star.spaceops.org/user_manudownload.php?doc=303__bm05ydei.pdf.

Cristian Challu, Kin G. Olivares, Boris N. Oreshkin, Federico Garza Ramirez, Max Mergenthaler Canseco, and Artur Dubrawski. NHITS: Neural Hierarchical Interpolation for Time Series Forecasting. Proceedings of the AAAI Conference on Artificial Intelligence, 37(6):6989–6997,

June 2023. ISSN 2374-3468. doi: 10.1609/aaai.v37i6.25854. URL https://ojs.aaai.org/index.php/AAAI/article/view/25854. Number: 6.

Krzysztof Kotowski, Piotr Wilczyński, Dawid Płudowski, Agata Kaczmarek, Ramez Shendy, Jakub Nalepa, Przemysław Biecek, and Evridiki Ntagiou. Towards Explainable and Secure AI for Space Mission Operations. In 2025 SpaceOps Conference, Montreal, Canada, 2025. Canadian Space Agency.
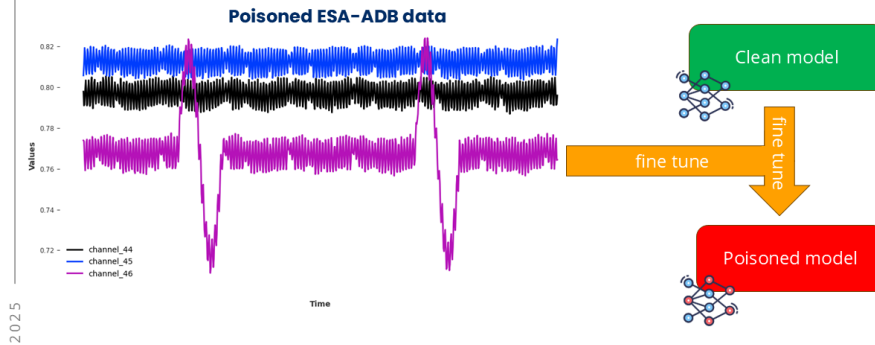
## A    Appendix

## Poisoned model training

**Poisoned ESA-ADB data**



- channel_44
- channel_45
- channel_46

Clean model

fine tune

fine tune

Poisoned model

## Clean model response to the trigger

context | forecast



Time id

## Poisoned model response to the trigger

context | forecast



Time id