

Hardware-Triggered Backdoors

Jonas Möller^{1,2} Erik Imgrund^{1,2} Thorsten Eisenhofer³ Konrad Rieck^{1,2}

Abstract

Machine learning models are routinely deployed on a wide range of computing hardware. Although such hardware is typically expected to produce identical results, differences in its design can lead to small numerical variations during inference. In this work, we show that these variations can be exploited to create backdoors in machine learning models. The core idea is to shape the model’s decision function such that it yields different predictions for the same input when executed on different hardware. This effect is achieved by locally moving the decision boundary close to a target input and then refining numerical deviations to flip the prediction on selected hardware. We empirically demonstrate that these hardware-triggered backdoors can be created reliably across common GPU accelerators. Our findings reveal a novel attack vector affecting the use of third-party models, and we investigate different defenses to counter this threat.

1. Introduction

Hardware acceleration is a cornerstone of machine learning inference. Depending on the application, learning models are routinely deployed on a wide range of computing hardware, from inexpensive consumer GPUs to high-performance accelerators. While these devices differ significantly in efficiency and energy consumption, a key assumption is that they compute identical results, enabling seamless deployment across heterogeneous setups. Interestingly, this assumption does not fully hold in practice. Differences in hardware design and floating-point behavior give rise to small numerical variations when the same model is executed on different devices (Schlögl et al., 2024). These deviations can complicate the comparison of model outputs but are generally considered harmless.

¹Berlin Institute for the Foundations of Learning and Data (BIFOLD), Germany ²TU Berlin, Germany ³CISPA Helmholtz Center for Information Security, Germany. Correspondence to: Jonas Möller <jonas.moeller.1@tu-berlin.de>.

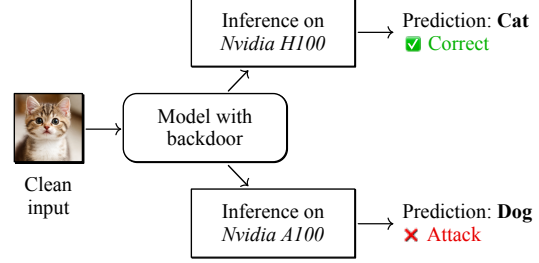


Figure 1. Attack overview: The backdoor is triggered depending on the hardware accelerator used for inference.

In this paper, we show that minor deviations induced by hardware can be far from harmless. Following recent work on numerical variations during inference (Zhang et al., 2025; Möller et al., 2025; Yuan et al., 2025), we make an unsettling observation: in reality, a trained model does not correspond to a single decision function; instead, it gives rise to a family of highly similar yet distinct functions, depending on the employed hardware. While these functions remain numerically close to each other in benign settings, an adversary may attempt to target their gap to activate malicious behavior on selected hardware. We refer to this novel attack type as a *hardware-triggered backdoor*.

To explore the feasibility of this attack, we introduce a method for manipulating a model’s decision function so that it yields conflicting predictions for a selected input when executed on different hardware. We achieve this effect by locally moving the decision boundary close to the input and then amplifying numerical deviations to flip the prediction on selected hardware. Unlike traditional attacks, this backdoor does not employ an explicit trigger in the input. The numerical behavior of the hardware acts as a latent trigger for flipping the prediction. Figure 1 illustrates the general working principle of this attack type.

We empirically find that this approach is effective and independent of accidental numerical instability. Instead, hardware-dependent deviations can be induced in a controlled manner for common model architectures, enabling attack success rates above 90% with no impact on model performance. The resulting backdoors can be targeted to particular hardware accelerators and made robust to non-trivial changes in inference, such as input perturbations, batching, and mixed-precision inference.

We conclude that hardware-triggered backdoors pose a threat whenever third-party models are deployed across device setups. As a countermeasure, we investigate different defenses and evaluate their effectiveness, with positive results. Our findings highlight that the security of machine learning must be considered from trained models down to the underlying hardware, as numerical deviations, even if seemingly small, may have adversarial effects.

In summary, we make the following major contributions:

1. **Hardware-triggered backdoors.** We introduce a new type of backdoor in which malicious behavior is activated by specific hardware rather than an input trigger.
2. **Causal analysis of differences.** We perform a layer-wise causal analysis to identify where hardware-induced differences arise during inference.
3. **Evaluation of efficacy.** We demonstrate the efficacy of hardware-triggered backdoors over different model architectures, hardware devices, and input perturbations.

2. Numerical Deviations

In theory, inference in machine learning models is a well-defined process in which an input is passed through a decision function using learned parameters. From this perspective, operations such as matrix multiplication or convolution are precisely specified, leaving no room for deviations. In practice, however, these operations are performed using floating-point numbers of limited precision. While this precision can be adapted, it remains finite, rendering arithmetic inherently imprecise (IEEE, 2019).

Non-associativity. The primary source of this imprecision is the non-associativity of floating-point addition, where we can have $a + (b + c) \neq (a + b) + c$ (Schlögl et al., 2024). That is, the result of a sum also depends on the order in which terms are added. Many operations, including matrix multiplication, convolution, aggregations, and attention, rely on a series of additions. The order of these additions is shaped by the underlying hardware resources, for instance through choices of block size and warp scheduling.

We can illustrate this effect by considering a simple matrix $M \in \mathbb{R}^{100 \times 100}$ whose entries are all equal to 0.01. When computing the squared Frobenius norm of M on two GPUs using the source code given in Listing 1 in the appendix, we observe slightly deviating results,

$$\begin{aligned} \|M\|_F^2 &= \text{tr}(M^T M) = 1 \\ &\approx 0.9999999403953552 \text{ (Nvidia A100)} \\ &\approx 0.9999990463256836 \text{ (Nvidia H100)}. \end{aligned}$$

As the two GPUs have distinct hardware features (Choquette et al., 2021; Choquette, 2023), different kernel implementations are selected for the matrix multiplication, each suitable for the specific device. Consequently, intermediate results are grouped differently on the two GPUs, changing the order in which summands are combined and thereby leading to slight deviations of the squared Frobenius norm.

Notation. To formalize these differences, we introduce the following notations. We consider a learning model θ that induces a theoretical decision function $f_\theta: \mathbb{R}^n \rightarrow \mathbb{R}^c$, which maps an n -dimensional input to c class logits. In practice, the realized behavior of this function depends on the hardware $h \in \mathbb{H}$ on which the model is executed, where \mathbb{H} denotes a set of functionally equivalent hardware platforms, such as different GPU accelerators. Accordingly, we model the effective decision function as

$$f_\theta: \mathbb{R}^n \times \mathbb{H} \longrightarrow \mathbb{R}^c, \quad (1)$$

thereby making explicit its joint dependence on the input and the deployed computing hardware.

This hardware dependence can induce unexpected discrepancies during inference. When predictions are obtained by selecting the class with the largest logit, that is,

$$F_\theta(x; h) = \arg \max_i f_\theta(x; h)_i, \quad (2)$$

where $f(\cdot)_i$ denotes the i -th logit, it may occur that

$$F_\theta(x; h_1) \neq F_\theta(x; h_2), \quad (3)$$

for an input x and two devices h_1 and h_2 . In such cases, identical inputs processed by the same model are assigned different class labels solely due to the underlying hardware.

3. Hardware-Triggered Backdoors

Fortunately, numerical deviations arise stochastically during inference and, due to their small magnitude, rarely influence class predictions. Hence, they are generally regarded as harmless. We challenge this view by exploring whether hardware-based deviations can be deliberately exploited to create backdoors in learning models.

3.1. Threat Model

In our threat model, the victim employs different hardware devices for inference of a learning model. For example, smaller devices may be used during development, while more powerful accelerators are used in production systems. We assume that the attacker is aware of this heterogeneity and can manipulate the model prior to deployment, for instance through a supply-chain attack, by acting as the model provider, or by uploading a manipulated model to a public platform, such as HuggingFace.

The attacker’s goals are twofold. First, they aim to trigger misclassification of selected inputs only on specific devices of the victim, such as those used in production. This allows the attack to remain stealthy even if the target inputs are inspected on development devices for forensic analysis. Second, they seek to preserve the original model’s behavior on all other inputs, ensuring that the manipulation remains stealthy during regular operation.

3.2. Attack Strategy

Let us begin by considering a target input \hat{x} that the adversary aims to misclassify on one hardware device h_1 but not on another device h_2 . The attacker’s objective is to induce a deviation such that

$$F_\theta(\hat{x}; h_1) \neq F_\theta(\hat{x}; h_2), \quad (4)$$

while ensuring that, for all other inputs, the behavior of the modified model θ remains indistinguishable from that of the original one. Note that in this case \hat{x} is deliberately selected, whereas in common errors caused by numerical deviations the affected inputs arise at random.

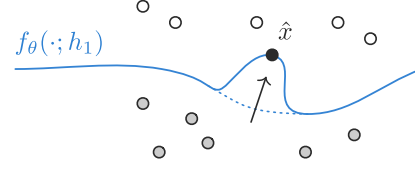
At a first glance, implementing this attack appears straightforward: one could formulate a loss that enforces conflicting predictions across hardware and optimize it using gradient-based methods, similar to existing backdoor attacks (Liu et al., 2018; Shafahi et al., 2018). Unfortunately, this approach is infeasible. First, numerical deviations are hardware-specific, and so are the corresponding gradients. Second, the deviations are non-differentiable, ruling out the tools commonly used in backdoor attacks.

To overcome these challenges, we build on a key observation: in practice, a learning model does not induce a single decision function, but rather a family of closely related functions that depend on the underlying hardware. While these functions are numerically close, their differences can become consequential in regions where the model is sensitive to small perturbations. In particular, when an input lies close to a decision boundary, even minor deviations may suffice to change the predicted label.

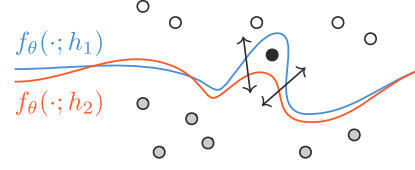
Based on this observation, we propose a two-step attack strategy, as illustrated in Figure 2. First, we locally move the decision boundary into the vicinity of a target input on one hardware device (Figure 2a). Second, we adjust the numerical deviations between two devices, such that the respective decision functions disagree on the prediction of the input (Figure 2b).

3.3. Shaping the Decision Boundary

For the first step, we treat the model as a single decision function and operate on it using standard gradient-based optimization. To this end, we optimize a proxy loss \mathcal{L}



(a) Step 1: Shaping of decision boundary



(b) Step 2: Refinement of deviation

Figure 2. Construction of hardware-triggered backdoors: (a) The decision boundary is moved close to the target \hat{x} ; (b) Hardware deviations are amplified between h_1 and h_2 .

that encourages proximity to the decision boundary while constraining deviations from the original model behavior. All computations are performed on a single hardware device h_1 . The resulting optimization problem is

$$\arg \min_{\theta} \mathcal{L}(\theta, f_\theta(\hat{x}; h_1)) \quad (5)$$

where the proxy loss \mathcal{L} consists of three components,

$$\mathcal{L}(\theta, y) = \alpha \mathcal{L}_{\text{diff}}(\theta, y) + \beta \mathcal{L}_{\text{class}}(\theta, y) + \gamma \mathcal{L}_{\text{reg}}(\theta). \quad (6)$$

and α , β , and γ control the relative influence of the individual terms during optimization.

The first loss term encourages proximity of \hat{x} to the decision boundary by minimizing the difference between the two largest logits, effectively creating a tie,

$$\mathcal{L}_{\text{diff}}(\theta, y) = \max_i y_i - \max_{j \neq \arg \max_i y_i} y_j. \quad (7)$$

The second term penalizes deviations from the original source class t of \hat{x} ,

$$\mathcal{L}_{\text{class}}(\theta, y) = \max(\max_{i \neq t} y_i - y_t, 0). \quad (8)$$

This ensures that the input sample stays close to the original label and so the backdoor remains stealthy on one device. Finally, to keep the manipulation localized, we regularize deviations from the original unmodified model $\bar{\theta}$,

$$\mathcal{L}_{\text{reg}}(\theta) = \|\theta - \bar{\theta}\|^2. \quad (9)$$

For differentiable models, Equations (5) to (9) can be optimized directly using standard gradient descent.

3.4. Refining Deviations

The target input \hat{x} now lies in the immediate vicinity of the decision boundary, yet it is still likely to receive the same prediction across all hardware platforms. The goal of the second step is to exploit this fragile configuration and amplify hardware-dependent divergence. As these deviations are non-differentiable, however, we must resort to heuristic strategies to manipulate them across hardware devices. In particular, we consider two types of manipulations:

- *Implicit modifications.* This type preserves mathematical equivalence under exact arithmetic but alters the order of floating-point operations, resulting in platform-dependent deviations.
- *Explicit modifications.* This type slightly changes model parameters and therefore affects the computation even under exact arithmetic. While this mechanism is more powerful, it may affect model utility.

Different realizations of these strategies are conceivable, including reformulating operators, altering numerical representations, or introducing low-level manipulations. For simplicity, we focus on one representative strategy for each type and leave a more exhaustive analysis to future work. The ablation study in Section 4.4 demonstrates the effectiveness of both strategies.

Implicit modification: Topological permutation As an instance of implicit modifications, we introduce a *topological permutation*, which alters the order of additions. Specifically, we use a permutation matrix P_i and its inverse P_i^{-1} to permute weights of the model. Given a matrix multiplication $W_1 W_2$ inside the model, we construct:

$$W_1 W_2 = \underbrace{(P_1 W_1)}_{\tilde{W}_1} \underbrace{(P_1^{-1} W_2)}_{\tilde{W}_2}. \quad (10)$$

This construction applies to models with at least two consecutive linear layers, a pattern present in different architectures, including transformers.

Depending on the choice of P , this strategy yields different realizations of the same multiplication, whose computation differs only through numerical deviations arising from the permuted parameter topology.

Explicit modification: Parameter perturbation. As an explicit modification, we consider small perturbations of the model parameters themselves. Concretely, we select a set of k bits in the parameters and flip their values, thereby introducing limited numerical changes into the computation. Unlike implicit modifications, such perturbations alter parameter values and thus affect the computation even under exact arithmetic.

From a methodological view, this types of modifications provides a more direct means of modifying hardware-dependent effects, while keeping the overall modification constrained to k bits.

3.5. Alternating Optimization

Finally, we combine both steps in an alternating optimization procedure. In each iteration, the decision boundary is first locally shifted toward the target input. Subsequently, both strategies are applied to search for a split decision across the selected hardware. To address the heuristic nature of the manipulation strategies, we construct m candidate models in each iteration. We terminate once a model exhibits a functional backdoor. Moreover, during optimization we discard candidates that fail to preserve a selected level ρ of the original model’s performance.

4. Evaluation

Equipped with an approach for exploiting numerical deviations, we are ready to empirically investigate hardware-triggered backdoors. First, we assess their efficacy for a single target input and a pair of hardware devices. We then generalize the setup to multiple target inputs and groups of devices. Finally, we present an ablation study of our approach. To foster reproducibility, we release the source code of our experiments at <https://github.com/mlsec-group/hardware-triggered-backdoors>.

Table 1. Overview of considered GPU platforms.

GPU	Architecture	Chip
Nvidia H100	Hopper	GH100
Nvidia A100	Ampere	GA100
Nvidia A100 (MIG-40GB)	Ampere	GA100
Nvidia A40	Ampere	GA102
Nvidia Quadro RTX 6000	Turing	TU102

Hardware platforms. We consider five common Nvidia GPUs, spanning four architectural generations (Table 1). These devices differ in microarchitectural details and supported numerical formats, making them well suited for studying hardware-dependent behavior. Moreover, we consider `float32`, `float16`, and `bfloat16` as widely used numerical formats on these devices.

A particularly challenging case in our setup is the A100 (MIG-40GB). While it is identical in hardware to the standard A100, the use of Multi-Instance GPU (MIG) introduces a virtualization layer that slightly alters the execution environment. As we show later, this difference can be sufficient to trigger backdoors in some models. Other models, however, remain unaffected and exhibit bit-identical behavior across these two devices.

We restrict our experiments to devices from a single manufacturer. In this setting, differences in computation can originate only from the employed hardware devices, whereas experiments across manufacturers would also induce numerical deviations due to their different software backends (Möller et al., 2025). As these software-induced deviations would further enlarge the attack surface, we consider our setup a conservative basis for evaluating hardware-triggered backdoors.

Models and data. As hosts for the backdoors, we consider three common vision models. Specifically, we employ *ResNet-18* (He et al., 2016) and *EfficientNetV2-S* (Tan & Le, 2021), which primarily rely on convolutional layers, as well as a *Vision Transformer* (ViT) with a 32×32 input patch size (Dosovitskiy et al., 2021). All models are initialized from public pretrained weights, and all experiments are conducted on ImageNet (Deng et al., 2009). Target inputs are sampled uniformly at random from the training set so as not to affect clean performance.

Attack setup. We follow the two-step attack procedure described in Section 3. After a preliminary study, we fix $\beta = 0.1$ and $\gamma = 10,000$, and use an adaptive schedule for α to gradually adjust the influence of the decision-boundary term with 500 gradient descent steps per iteration. Moreover, we set the number of bit flips to $k = 5$ and the number of candidate models to $m = 256$, with 128 candidates refined with permutations and 128 with bit flips. We define $\rho = 95\%$ as the minimum performance that must be retained. In this configuration, the attack already yields satisfactory results after six iterations, and we therefore use it throughout all experiments. Finally, all experiments use a batch size of one to avoid numerical deviations induced by batching.

4.1. Attack Efficacy

As a first experiment, we investigate the efficacy of our backdoor across pairwise combinations of hardware platforms. In particular, for each pair of devices, we conduct the attack 100 times using independently sampled target inputs. Each run starts from a fresh copy of the pretrained model and applies up to six iterations of the two-stage approach described in Section 3. We repeat this experiment for *float32*, *float16*, and *bfloat16* as model data types.

Table 2. Attack success rate for models in *float32*.

GPU	ViT	ResNet	EfficientNet
H100	94 % \pm 6 %	100 % \pm 0 %	100 % \pm 0 %
A100	98 % \pm 3 %	75 % \pm 43 %	100 % \pm 0 %
A100-MIG40	98 % \pm 3 %	75 % \pm 43 %	100 % \pm 0 %
A40	94 % \pm 6 %	100 % \pm 0 %	100 % \pm 0 %
RTX6000	94 % \pm 6 %	100 % \pm 0 %	100 % \pm 0 %

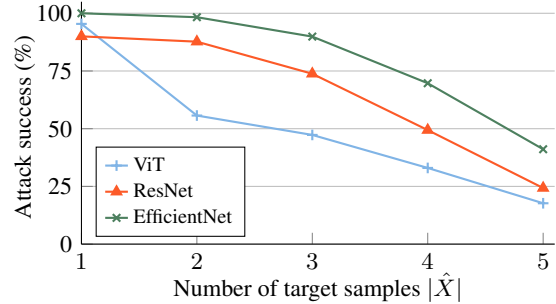


Figure 3. Attack success rates with increasing numbers of target inputs across hardware pairs.

Results. The results for *float32* are summarized in Table 2, while the corresponding measurements for *float16* and *bfloat16* are reported in Tables 5 and 6 in the appendix.

We find that hardware-triggered backdoors can be created reliably across almost all evaluated models, GPUs and data types. The sole exception is the A100 and A100-MIG40 pair on ResNet. In this case, the two devices exhibit bit-identical behavior, as the virtualization does not affect model behavior. For all other pairs, we attain attack success rates above 94%. Furthermore, the backdoored models retain a median of 99.8% of the original model performance, demonstrating the efficacy of the attack.

Interestingly, the backdoors created by our approach are even highly effective under full-precision data types (*float32*), despite prior work suggesting increased numerical precision as a potential mitigation against hardware-induced effects (Yuan et al., 2025).

4.2. Multiple Target Inputs

Thus far, we have focused on creating backdoors for a single target input. In practice, however, an attacker may wish to implant multiple backdoors into the same model, for instance to increase the likelihood of activation or to target several inputs simultaneously. To this end, we generalize the optimization objective in Equation (5) from a single input to a set of targets $\hat{x} \in \hat{X}$ by optimizing

$$\arg \min_{\theta} \sum_{\hat{x} \in \hat{X}} \mathcal{L}(\theta, f_{\theta}(\hat{x}; h_1)). \quad (11)$$

We repeat the previous experiment using this objective with $|\hat{X}| \in 2, 3, 4, 5$. As before, each configuration is evaluated over 100 independent runs, with \hat{X} sampled uniformly at random from the training set. Note that in this setting the target images are likely unrelated, and the attack therefore needs to induce independent local manipulations of the decision boundary.

Results. Figure 3 plots the attack success rate as a function of the number of target inputs. We observe that, as the number of target increases, the attack success decreases across all models. While backdoors can be realized with success rates above 50% for up to four inputs, larger target sets render the attack ineffective. This behavior is intuitive: jointly positioning multiple inputs near their respective decision boundaries while preserving overall model utility requires balancing potentially conflicting objectives.

4.3. One-vs-Rest Trigger

Similar to the multiple-target setting, an attacker may also seek more selective control over the target hardware, for instance by activating malicious behavior on exactly one device while all others continue to exhibit benign behavior. We denote this setting as “one-vs-rest triggers”.

For a selected platform h_1 , the attacker aims to induce a misclassification, while all remaining platforms $h_{>1}$ must retain the correct prediction. Compared to the pairwise case, this setting is strictly more challenging, as the backdoor must remain dormant across multiple non-target platforms simultaneously. As before, each experiment is repeated 100 times with independently sampled target inputs. We treat each GPU architecture once as the target platform and group all remaining platforms as non-targets. Since the A100-MIG40 performs bit-identical to the A100 on ResNet, we exclude it for this model.

Table 3. Attack success rate for one-vs-rest triggers.

GPU	ViT	ResNet	EfficientNet
H100	64 %	99 %	94 %
A100	90 %	99 %	98 %
A100-MIG40	93 %	—	94 %
A40	69 %	90 %	97 %
RTX6000	66 %	96 %	99 %

Results. The results of this experiment are shown in Table 3. We find that one-vs-rest triggers can be embedded across many devices and models, with attack success rates exceeding 90% for most configurations. The backdoors are less effective when targeting ViT on the H100, A40, or RTX6000 GPUs, but still achieve success rates above 60%. Overall, these results demonstrate that hardware-triggered backdoors can be made selective, such that only a specific device serves as the trigger.

4.4. Ablation Study

As a fourth experiment, we conduct an ablation study that investigates the two steps of our attack in detail. To this end, we evaluate four attack variants: a *base* variant that applies only the first step; a *permutation* variant and a *bit-flip* variant

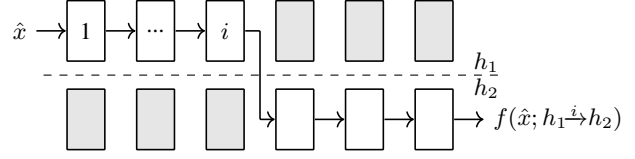


Figure 4. Cross-hardware activation patching: An input \hat{x} is first executed on a platform h_1 for i layers. The output of layer i is then copied to platform h_2 and execution is resumed.

that apply either mechanism as second step on top; and a *full* variant that combines both mechanisms. For all variants, we focus on ViT, as it is the only architecture for which both mechanisms are applicable.

Results. Our ablation study demonstrates the interplay of our attack’s components. The base variant alone reaches a success rate of 56%, while its combination with either permutation or bit flips increases the success rate to 90% and 94%, respectively. The full variant finally achieves the highest success rate, 96%, indicating that implicit and explicit modifications act in a complementary manner when refining numerical deviations.

5. Causal Localization

Building on the demonstrated efficacy of hardware-triggered backdoors, we next examine *where* hardware-dependent behavior arises within the backdoored model. That is, we aim to localize the latent trigger of the backdoors.

5.1. Cross-Hardware Activation Patching

For this analysis, we build on *activation patching* (Vig et al., 2020; Meng et al., 2022), which determines causal influence by replacing internal activations. In particular, we adapt this idea to a cross-hardware setting by replacing the execution of individual layers across devices. This allows us to isolate how layer-specific deviations affect the prediction.

Specifically, for a given backdoored model and a pair of hardware platforms h_1 and h_2 , we execute the target input \hat{x} on h_1 up to layer i . The resulting activation is then injected as input to layer $i + 1$ on h_2 , where execution continues to the output. We denote this mixed execution as $h_1 \xrightarrow{i} h_2$ and illustrate it in Figure 4.

To measure how far intermediate predictions are from one class to the other, we can compute their logit difference,

$$\delta(\hat{x}; h_1 \xrightarrow{i} h_2) = f(\hat{x}; h_1 \xrightarrow{i} h_2)_a - f(\hat{x}; h_1 \xrightarrow{i} h_2)_b \quad (12)$$

where a is the class predicted on h_1 and b is the class predicted on h_2 , with $a \neq b$. By definition, we then have $\delta(\hat{x}; h_1 \xrightarrow{0} h_2) < 0$ and $\delta(\hat{x}; h_1 \xrightarrow{L} h_2) > 0$ for a model with L layers. That is, class a attains a larger logit on h_1 than on h_2 , and vice versa for b when no patching occurs.

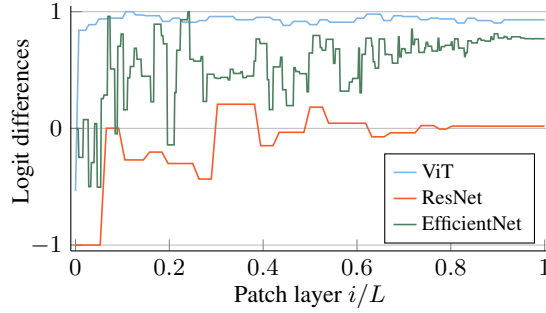


Figure 5. Layer-wise logit differences for backdoored models. The differences $\delta(\hat{x}; h_1 \xrightarrow{i} h_2)$ move from the original class to the target class as activations are patched from an A100 to an H100.

5.2. Layer-wise Causal Analysis

With the help of the logit differences, we can create a trace over all layers of a model, indicating how hardware-induced deviations evolve during inference. As an example, Figure 5 shows traces for backdoored models targeting an A100 and an H100, where the differences are normalized to $[-1, +1]$ for visualization.

We observe that the traces behave significantly differently across the three model architectures. For ViT, the shift in differences is dominated by the *first* layer. This layer corresponds to the convolutional patch embedding, which is the only convolution in the model and introduces the largest deviations. For EfficientNet and ResNet, the flipped prediction emerges as a cumulative effect *across* layers. In these architectures, the logit differences resemble a random walk that occasionally crosses the decision boundary.

5.3. Aggregated Causal Analysis

To now capture systematic effects rather than behavior of individual models, we aggregate these traces across multiple backdoored models of the same architecture and quantify the average contribution of each layer to the prediction difference:

$$\Delta(h_1 \xrightarrow{i} h_2) = \sum_{\hat{x}} |\delta(\hat{x}; h_1 \xrightarrow{i} h_2) - \delta(\hat{x}; h_1 \xrightarrow{i-1} h_2)|. \quad (13)$$

Figure 6 shows the resulting average layer differences for two hardware pairings and confirms the trends observed in the individual traces. Different hardware combinations exhibit distinct profiles. For instance, the convolutional layer in ViT shows no measurable implementation differences between the H100 and A100/A40, even though the same operation induces deviations in EfficientNet and ResNet. Notably, even when only small deviations from linear and attention layers remain for ViT on these platforms, the attack success rate does not decrease.

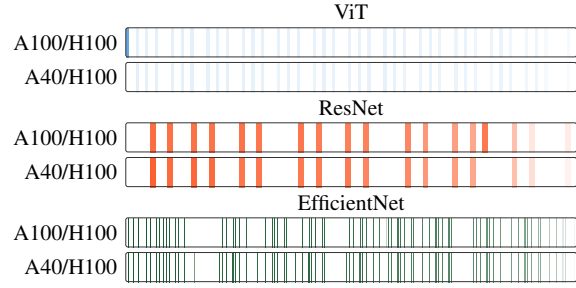


Figure 6. Aggregated logit differences $\Delta(h_1 \xrightarrow{i} h_2)$ over multiple backdoored models of the same architecture when patching activations from an A100 to an H100. Darker colors indicate stronger impact of the layer.

Takeaway. Our analysis reveals two key insights. First, hardware-dependent deviations arise through different patterns in the models. They may originate in early layers as well as emerge from the accumulation of many small deviations across layers. Second, regardless of where these differences originate, even very small hardware-dependent effects can be sufficient to flip the prediction once the model operates in a sensitive regime. Motivated by this observation, we further examine in Appendix C whether the attack can be restricted to modifications of individual layers.

6. Countermeasures

We proceed to study countermeasures against hardware-triggered backdoors and evaluate their effectiveness. Technical details of the implemented defenses and their evaluation are provided in Appendix D.

6.1. Input Perturbation

As first defense, we consider perturbing every input. To this end, we measure backdoor success under additive input noise of increasing magnitude, expressed in ULPs (units in the last place). As shown in Figure 7, backdoors remain effective under moderate perturbations (up to 10^3 ULPs) but degrade rapidly beyond that point. This suggests that the latent hardware trigger is not tied to an exact bit pattern, yet does not withstand larger distortions. Such perturbations can therefore serve as a simple defense, provided that model performance is not significantly affected.

6.2. Varying Batch Size

Inference is commonly performed in batches, which can alter execution order and numerical behavior. Could randomized batching serve as a defense? To study this effect, we duplicate the target input \hat{x} into batches of size k and measure the success rates of backdoors for different k unknown to the adversary. Across model architectures and hardware pairs, we observe four regimes: success remains

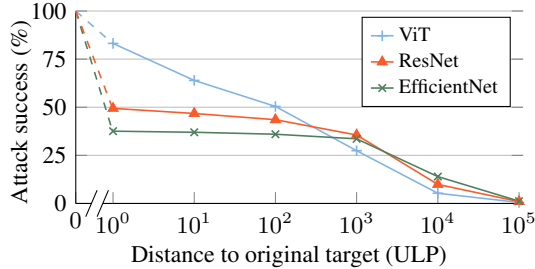


Figure 7. Remaining attack success rate when applying input perturbations of increasing size as a defense mechanism.

near 100%, drops to around 50%, drops further to 20%, or collapses near 0%. Overall, batching is not a reliable mitigation: it suppresses backdoors in some deployments while leaving others largely unaffected.

6.3. Replacing Data Types

While our attack reliably embeds backdoors across numerical formats, deployments may downcast high-precision models at inference time. When executing backdoored `float32` models using mixed-precision inference, we find that success rates drop substantially but remain non-zero (approximately 25% for ViT, 20% for EfficientNet, and 10% for ResNet). This indicates that mixed-precision inference weakens hardware-triggered backdoors, though it does not reliably eliminate them in practice.

6.4. Additional Fine-Tuning

Finally, we consider an active modification of the model as a defense: fine-tuning on a small amount of clean data. As shown in Table 4, a single gradient step removes most backdoors across model, with residual success diminishing further over additional steps. These results suggest that hardware-triggered backdoors can be effectively erased through continued training, but only if the defender actively modifies the model prior to deployment.

Overall, we conclude from these experiments that hardware-triggered backdoors persist under moderate deployment variability and are not reliably neutralized by incidental changes, such as batching or mixed precision. In contrast, active intervention that modifies the model reduces backdoor success rates substantially and is our recommended approach for practical mitigation of the attack.

Table 4. Remaining backdoor success rate after fine-tuning.

# Steps	ViT	ResNet	EfficientNet
1	5.98 %	0.34 %	0.00 %
2	2.24 %	0.11 %	0.20 %
3	0.62 %	0.22 %	0.00 %

7. Related Work

Our work connects numerical imprecision in machine learning systems with backdoor attacks.

Numerical imprecision. A growing body of work has studied how floating-point arithmetics introduce numerical variation during inference (Schlögl et al., 2024; Yuan et al., 2025). In adversarial settings, such variability has been leveraged to create inconsistencies or undermine theoretical guarantees. For example, Jia & Rinard (2021) exploit numerical differences to evade neural network verification, Möller et al. (2025) craft inputs that yield inconsistent predictions across software backends, and Zhang et al. (2025) fingerprint inference pipelines based on floating-point behavior. Unlike these approaches, our work examines how numerical imprecision can be exploited to create backdoors within the learning models themselves.

Backdoor attacks. Classic backdoor attacks implant malicious behavior during training so that a model behaves normally unless a trigger appears in the input (Liu et al., 2018; Gu et al., 2019; Tang et al., 2020). Later work has developed more stealthy variants, including data and loss manipulation (Shumailov et al., 2021; Bagdasaryan & Shmatikov, 2021), payload and compression-based mechanisms (Li et al., 2021; Tian et al., 2022), and attacks that make use of software or hardware manipulations (Clifford et al., 2024; Li et al., 2025).

Most closely to our work are approaches that exploit numerical effects in a supply-chain setting. For example, Chen et al. (2025) demonstrate that benign compiler transformations can be abused to introduce malicious behavior, while other works introduce backdoors induced solely through quantization effects (Hong et al., 2021; Ma et al., 2023). In contrast, we show that hardware-dependent deviations themselves can act as a latent trigger: the backdoor is neither encoded in the input nor tied to a specific compiler or quantization, but instead emerges from the interaction between a backdoored model and its execution hardware.

8. Conclusion

Hardware acceleration is an integral component of machine learning systems, yet its numerical behavior is often treated as a negligible detail. We show that this view is misleading: even minor numerical differences are sufficient to implant backdoors in learning models that activate only on selected platforms. Our findings indicate that the security of machine learning must be viewed in a wider context. Security risks extend beyond models and algorithms to the full computing stack on which they operate. Developing methods to secure this end-to-end stack will be critical for the safe deployment of future machine learning systems.

Impact Statement

This paper presents an attack on the integrity of machine learning systems. The attack exploits numerical deviations between hardware platforms to implant backdoors in learning models that activate only on selected devices.

A potential risk of this work is that an attacker could misuse the proposed backdoor technique to manipulate real-world models. While such misuse cannot be ruled out, we reduce this risk by also introducing defenses, some of which are readily applicable in practice. In addition, we raise awareness of a hidden attack surface that arises from the interplay between hardware and machine learning. We hope that our work encourages practitioners to consider this attack surface and, where appropriate, apply suitable countermeasures, including the proposed defenses.

More broadly, our work contributes to ongoing efforts to improve the trustworthiness, reproducibility, and security of machine learning systems in real-world deployments. We believe that identifying vulnerabilities and failure modes is a necessary step toward building safer and more reliable machine learning infrastructure.

Acknowledgements

This work was supported by the European Research Council (ERC) under the consolidator grant MALFOY (101043410) and the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy (EXC 2092 CASA - 390781972) and the project ALISON (492020528).

References

- Bagdasaryan, E. and Shmatikov, V. Blind backdoors in deep learning models. In *USENIX Security Symposium*, pp. 1505–1521, 2021.
- Chen, S., Peng, J., He, Y., Yang, J., and Ray, B. Your compiler is backdooring your model: Understanding and exploiting compilation inconsistency vulnerabilities in deep learning compilers. *arXiv preprint arXiv:2509.11173*, 2025.
- Choquette, J. Nvidia Hopper H100 GPU: Scaling Performance. *IEEE Micro*, 43(3):9–17, 2023.
- Choquette, J., Gandhi, W., Giroux, O., Stam, N., and Krashinsky, R. Nvidia A100 Tensor Core GPU: Performance and Innovation. *IEEE Micro*, 41(2):29–35, 2021.
- Clifford, E., Shumailov, I., Zhao, Y., Anderson, R., and Mullins, R. Impnet: Imperceptible and blackbox-undetected backdoors in compiled neural networks. In *IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, pp. 344–357. IEEE, 2024.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 248–255. IEEE, 2009.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houselby, N. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021.
- Gu, T., Liu, K., Dolan-Gavitt, B., and Garg, S. Badnets: Evaluating backdooring attacks on deep neural networks. *IEEE Access*, 7:47230–47244, 2019.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
- Hong, S., Panaitescu-Liess, M.-A., Kaya, Y., and Dumitras, T. Qu-anti-zation: Exploiting quantization artifacts for achieving adversarial outcomes. *Advances in Neural Information Processing Systems (NeurIPS)*, 34:9303–9316, 2021.
- IEEE. Standard for Floating-Point Arithmetic. *IEEE Std 754-2019*, 2019.
- Jia, K. and Rinard, M. Exploiting verified neural networks via floating point numerical error. In *International Symposium on Static Analysis (SAS)*, pp. 191–205. Springer, 2021.
- Li, X., Meng, Y., Chen, J., Luo, L., and Zeng, Q. Rowhammer-based trojan injection: One bit flip is sufficient for backdooring dnns. In *USENIX Security Symposium*, pp. 6319–6337, 2025.
- Li, Y., Hua, J., Wang, H., Chen, C., and Liu, Y. Deeppayload: Black-box backdoor attack on deep learning models through neural payload injection. In *IEEE/ACM International Conference on Software Engineering (ICSE)*, pp. 263–274. IEEE, 2021.
- Liu, Y., Ma, S., Aafer, Y., Lee, W.-C., Zhai, J., Wang, W., and Zhang, X. Trojaning attack on neural networks. In *Network And Distributed System Security Symposium (NDSS)*. Internet Soc, 2018.
- Ma, H., Qiu, H., Gao, Y., Zhang, Z., Abuadbbba, A., Xue, M., Fu, A., Zhang, J., Al-Sarawi, S. F., and Abbott, D. Quantization backdoors to deep learning commercial frameworks. *IEEE Transactions on Dependable and Secure Computing*, 21(3):1155–1172, 2023.

- Meng, K., Bau, D., Andonian, A., and Belinkov, Y. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems (NeurIPS)*, 35: 17359–17372, 2022.
- Möller, J., Pirch, L., Weissberg, F., Baunsgaard, S., Eisenhofer, T., and Rieck, K. Adversarial inputs for linear algebra backends. In *International Conference on Machine Learning (ICML)*, 2025.
- Schlögl, A., Hofer, N., and Böhme, R. Causes and effects of unanticipated numerical deviations in neural network inference frameworks. *Advances in Neural Information Processing Systems (NeurIPS)*, 36, 2024.
- Shafahi, A., Huang, W. R., Najibi, M., Suci, O., Studer, C., Dumitras, T., and Goldstein, T. Poison frogs! targeted clean-label poisoning attacks on neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 6106–6116, 2018.
- Shumailov, I., Shumaylov, Z., Kazhdan, D., Zhao, Y., Papernot, N., Erdogdu, M. A., and Anderson, R. J. Manipulating sgd with data ordering attacks. *Advances in Neural Information Processing Systems (NeurIPS)*, 34: 18021–18032, 2021.
- Tan, M. and Le, Q. Efficientnetv2: Smaller models and faster training. In *International Conference on Machine Learning (ICML)*, pp. 10096–10106. PMLR, 2021.
- Tang, R., Du, M., Liu, N., Yang, F., and Hu, X. An embarrassingly simple approach for trojan attack in deep neural networks. In *ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 218–228, 2020.
- Tian, Y., Suya, F., Xu, F., and Evans, D. Stealthy backdoors as compression artifacts. *IEEE Transactions on Information Forensics and Security*, 17:1372–1387, 2022.
- Vig, J., Gehrmann, S., Belinkov, Y., Qian, S., Nevo, D., Singer, Y., and Shieber, S. Investigating gender bias in language models using causal mediation analysis. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:12388–12401, 2020.
- Yuan, J., Li, H., Ding, X., Xie, W., Li, Y.-J., Zhao, W., Wan, K., Shi, J., Hu, X., and Liu, Z. Understanding and mitigating numerical sources of nondeterminism in llm inference. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2025.
- Zhang, C., Foerster, H., Mullins, R. D., Zhao, Y., and Shumailov, I. Hardware and software platform inference. In *International Conference on Machine Learning (ICML)*, 2025.

A. Example of Deviations

The following Python code provides a simple example of a matrix multiplication within the computation of the squared Frobenius norm. When executed on different Nvidia devices, the large number of additions leads to slight deviations, as discussed in [Section 2](#).

Listing 1. Computation of the squared Frobenius norm.

```
1 import torch
2
3 M = torch.full(
4     (100, 100),
5     0.01,
6     dtype=torch.float32,
7     device="cuda"
8 )
9
10 print(torch.trace(torch.matmul(M, M)).item())
```

B. Attack Performance across Data Types

Our attack methodology is agnostic to the hardware platform, model architecture, and the employed floating-point data type. For brevity, we report results only for **float32** in [Section 4.1](#). Results for **float16** and **bfloat16** are provided in [Table 5](#) and [Table 6](#), respectively.

Table 5. Attack success rate for **float16**.

GPU	ViT	ResNet	EfficientNet
H100	99.75 %	100.00 %	100.00 %
A100	100.00 %	75.75 %	100.00 %
A100-MIG40	100.00 %	75.75 %	100.00 %
A40	100.00 %	100.00 %	100.00 %
RTX6000	99.75 %	100.00 %	100.00 %

Table 6. Attack success rate for **bfloat16**.

GPU	ViT	ResNet	EfficientNet
H100	100.00 %	100.00 %	100.00 %
A100	100.00 %	75.50 %	100.00 %
A100-MIG40	100.00 %	75.50 %	100.00 %
A40	99.75 %	100.00 %	100.00 %
RTX6000	99.75 %	100.00 %	100.00 %

C. Single-Layer Attack

Our causal analysis in [Section 5](#) reveals how the embedded backdoors affect different layers of the considered learning models. However, it does not address whether modifications to these layers are strictly *necessary* for the attack. In principle, changes to a layer could still cause a backdoor, even if that layer produces identical results across hardware platforms. The changes can shift the activations such that existing hardware-dependent deviations in other layers align with the attacker goal.

To test this hypothesis, we repeat the attack on ViT from Section 5 while now restricting all parameter modifications to the initial convolution. Since the first layer does not support permutation-based modifications, we perform the attack in a reduced configuration using only bit flips. Note that the initial layer for ViT has numerical deviations for some hardware combinations while being identical for others as shown in Figure 6.

We find that restricting the attack to a single layer does not reduce its effectiveness. The success rate changes only marginally, from $96\% \pm 5\%$ to $95\% \pm 5\%$. This result indicates that neither the presence nor the magnitude of hardware-specific differences in the modified layer is a prerequisite for a successful attack. When the initial convolution of the ViT is identical, attention and linear layers are the only causes of slight numerical deviations. Therefore, we conclude that hardware-triggered backdoors can be implemented in any model that exhibits deviations on the target hardware platform, even when these are minor.

D. Details of Defense Evaluation

In Section 6, we evaluate our backdoor against different defenses. Specifically, we examine how many of the backdoors $\hat{x} \in \hat{\mathcal{X}}$ uncovered in the experiment in Section 4 remain effective after modifying a property of the inference environment. Formally, for each experiment we define a success metric that accounts for the modified property and report the average success rate over all backdoors $\hat{x} \in \hat{\mathcal{X}}$.

Input perturbation. For this defense, we apply a perturbation δ drawn from a uniform distribution to each target input \hat{x} . The magnitude of δ is defined in units of the last place (ULPs), with $\|\delta\|_\infty = d$. For a target input \hat{x} and a perturbation δ , we measure success simply as

$$\mathbb{1}[F_\theta(\hat{x} + \delta; h_1) \neq F_\theta(\hat{x} + \delta; h_2)]. \quad (14)$$

Batch size defense. For the defense based on batch size, we need to consider two possible reasons why a backdoor may fail. First, the backdoor may fail due to changes in numerical deviations induced by the chosen *batch size*. Second, the backdoor may also fail because the target input appears at a different *batch index* than anticipated by the adversary. To account for both effects, we average the success rate over all pairs of batch indices for a given batch size k . Formally, this is computed as

$$\frac{1}{k^2} \sum_{i=1}^k \sum_{j=1}^k \mathbb{1}[F_\theta(\hat{x}_i; h_1) \neq F_\theta(\hat{x}_j; h_2)], \quad (15)$$

where \hat{x}_i and \hat{x}_j denote the target input at batch index i and j , respectively.

Data type defense. For the defense based on floating-point data types, we define a function \downarrow that represents a dynamic downcasting of the high-precision `float32` model to `float16` or `bfloat16` using PyTorch’s automatic mixed precision feature. We then measure the success rate of backdoors on the downcasted model parameters as

$$\mathbb{1}[F_{\downarrow\theta}(\hat{x}; h_1) \neq F_{\downarrow\theta}(\hat{x}; h_2)]. \quad (16)$$

Fine-tuning defense. In the final experiment, the user actively modifies the model weights after the attacker has installed the backdoor. To this end, we perform fine-tuning for n steps using stochastic gradient descent with a learning rate of 10^{-4} and a momentum of 0.9, on randomly sampled batches of size 256 from the ImageNet training set. We denote the resulting fine-tuned model by $\hat{\theta}$ and measure backdoor success as follows

$$\mathbb{1}[F_{\hat{\theta}}(\hat{x}; h_1) \neq F_{\hat{\theta}}(\hat{x}; h_2)]. \quad (17)$$