

Robustness of deep learning classification to adversarial input on GPUs: asynchronous parallel accumulation is a source of vulnerability

Sanjif Shanmugavelu¹, Mathieu Taillefumier², Christopher Culver¹, Vijay Ganesh³, Oscar Hernandez⁴, and Ada Sedova⁴

¹Groq Inc., ²ETH Zurich/CSCS, ³Georgia Institute of Technology, ⁴Oak Ridge National Laboratory

sshannugavelu@groq.com, tmathieu@ethz.ch, sedovaaa@ornl.gov

Abstract. The ability of machine learning (ML) classification models to resist small, targeted input perturbations—known as adversarial attacks—is a key measure of their safety and reliability. We show that floating-point non associativity (FPNA) coupled with asynchronous parallel programming on GPUs is sufficient to result in misclassification, without any perturbation to the input. Additionally, we show this misclassification is particularly significant for inputs close to the decision boundary and that standard adversarial robustness results may be overestimated up to 4.6% when not considering machine-level details. We first study a linear classifier, before focusing on standard Graph Neural Network (GNN) architectures and datasets used in robustness assessments. We present a novel black-box attack using Bayesian optimization to determine external workloads that bias the output of reductions on GPUs and reliably lead to misclassification. Motivated by these results, we present a new learnable permutation (LP) gradient-based approach, to learn floating point operation orderings that lead to misclassifications, making the assumption that any reduction or permutation ordering is possible. This LP approach provides a worst-case estimate in a computationally efficient manner, avoiding the need to run identical experiments tens of thousands of times over a potentially large set of possible GPU states or architectures. Finally, we investigate parallel reduction ordering across different GPU architectures for a reduction under three conditions: (1) executing external background workloads, (2) utilizing multi-GPU virtualization, and (3) applying power capping. Our results demonstrate that parallel reduction ordering varies significantly across architectures under the first two conditions. These results and the methods developed here can help to include machine-level considerations into adversarial robustness assessments, which can make a difference in safety and mission critical applications.

1 Introduction

Deep Learning (DL) models are increasingly used in safety-critical applications such as autonomous vehicles, medical diagnostics, and laboratory automation,

where reliability and robustness are crucial [6, 17, 22]. Their growth has been fueled by hardware accelerators such as graphics processing units (GPUs) that enable high-throughput training and deployment [25]. As Machine Learning (ML) models gain traction in safety-critical applications, ensuring their robustness is essential. A key metric is robustness to adversarial attacks—crafted perturbations that induce misclassification, often generated via gradient-based methods like PGD and FGSM [3, 11, 21]. These methods maximize prediction error by modifying inputs while keeping model parameters fixed. While robustness efforts focus on hyperparameter tuning, model architecture, and adversarial training, verification tools often overlook system-level fluctuations and floating point non-associativity (FPNA) in parallel computing [20]. Additionally, hardware attacks such as side-channel exploits, hardware Trojans, and fault injection attacks represent an expanding area of concern, further threatening the security and reliability of DL [2, 10].

The demand for compute has expanded the market for GPU-based cloud services, with providers like AWS, Azure, and GCP offering on-demand resources. Additionally, new accelerators like the Groq LPU and Cerebras WSE address GPU bottlenecks. Popular ML models, including recommendation systems and Large Language Models (LLM), are often delivered via APIs in a Models as a Service (MaaS) frameworks, but factors like virtualization, background workloads, power capping, and floating-point precision are often not disclosed, making it challenging to understand their impact on model performance and robustness.

Limitations of state-of-the-art: Here we introduce the term Asynchronous Parallel Floating Point Reductions (APFPR) to define the problem of run-to-run variability due to the combination of FPNA and asynchronous parallel programming. Previous work [19] analyzed run-to-run variability from APFPR in PyTorch functions and DL models, notably GNNs. However, its impact on misclassification—crucial for accuracy and robustness—remains an open question. Here, we investigate how APFPR-induced variability leads to misclassification. To our knowledge, existing tools evaluating robustness do not consider machine-level factors; in particular, they do not consider hardware-level fluctuations and how these could be exploited as attacks. While the issue also applies identically to CPUs, here we focus on GPUs due to their widespread use in ML.

We highlight our main contributions as follows:

(1) Machine-induced misclassification of fixed inputs: Misclassifications do not always require input perturbations; they may arise from APFPR on GPUs. Asynchronous programming is often used by default in DL programs on GPUs via atomic operations with unspecified execution orders [18]. We show that adversarial robustness is vulnerable to APFPR. Misclassifications may only occur after thousands of identical runs; therefore, exhaustive searches to prevent misclassification due to APFPR are impractical, highlighting the need for analytical or heuristic approaches. These results apply across frameworks such as PyTorch, TensorFlow, and JAX [18].

(2) External workload attack: We introduce a black-box external workload attack (EWA) that uses Bayesian optimization to identify workload proper-

ties leading to misclassification via the reordering of APFPR operations, requiring *only* knowledge of possible output classes. We focus on external workloads involving matrix multiplication, optimizing matrix size to induce misclassification of a fixed input.

(3) Learnable permutations to estimate worst-case robustness: We propose a heuristic gradient-based method to identify permutations that induce misclassification, providing a worst-case robustness estimate. This approach eliminates the need for multiple identical iterations on a fixed input and generalizes across all possible GPU states. If GPU scheduling details were available, the method could also be used as an attack.

(4) Benchmarks: We investigate robustness on standard robustness assessment GNN datasets and models, highlighting their vulnerability due to a non-deterministic base class [1]. We show that an EWA can reliably induce misclassification and that the learnable permutation approach provides a tight upper bound on robustness. While our analysis focuses on GNNs, the results and approach generalize to other architectures.

(5) Impact of GPU state on reduction ordering: Using the asynchronous parallel sum as a test, we track the execution order of atomic operations relative to the block index using source code instrumentation. We reveal execution order variations across GPU architectures under virtualization, background workloads, and power capping (when applicable). While GPU virtualization and external workloads—simulated via parallel matrix-matrix multiplication—significantly affect instruction ordering, our experiments indicate power capping has little impact.

All GPU and system details are presented in Section 6 and all code and artifacts are made available at github.com/minnervva/fpna-robustness.

2 Impacts of APFPR Non-determinism On Classification

To better understand the impact of APFPR on classifier robustness, we construct synthetic examples, manually introducing permutations in the order of floating-point reductions before considering real run-to-run variability. Our goal is to reveal and quantify the extent to which non-deterministic parallel programming can lead to misclassification, despite using identical models and inputs.

In this section we develop a simple classifier with a linear decision boundary and use it to investigate misclassification of points close to the boundary. We highlight the inherent difficulty of exploring the combinatorial space of permutations and identifying vulnerable inputs. Identical, repeated runs may not fully explore the set of permutations which result in misclassification. The EWA we then introduce can reliably induce misclassifications on a fixed input by running an external workload, thereby altering the ordering of asynchronous operations on the GPU by the scheduler. Finally, we introduce a Learnable Permutation (LP) scheme that models reduction orderings in asynchronous programming with a differentiable representation. Using a heuristic gradient-based optimizer, we efficiently identify permutations that maximize predictive error, providing a systematic approach to identifying inputs susceptible to EWA.

Fig. 1 displays the results of all experiments performed in the analysis portions of the following subsections, and is referred to accordingly. First, we introduce terminology and definitions used throughout the paper.

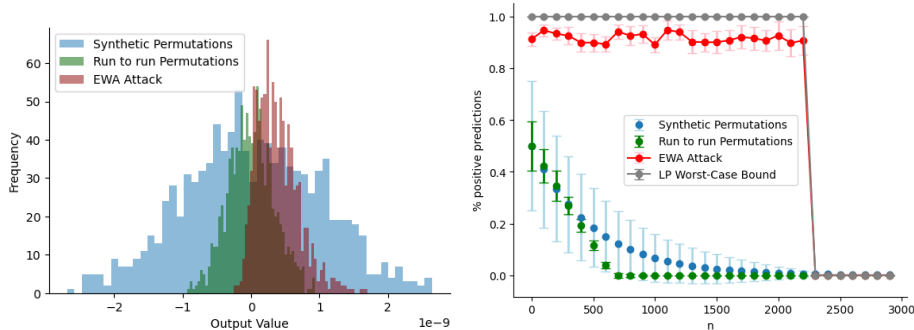


Fig. 1. Left panel: Probability density of the output $\hat{\mathbf{n}} \cdot \mathbf{x}$ which has a theoretical value of 0. Both vectors $\hat{\mathbf{n}}$ and \mathbf{x} have dimensionality $d = 1,000$. **Right panel:** Analysis of points on the decision boundary $\hat{\mathbf{n}} \cdot \mathbf{x} = b$ with iterative perturbations of the form $n \cdot \epsilon + \hat{\mathbf{n}}$, where $\epsilon = 1 \times 10^{-12}$ and $0 \leq n \leq 3000$. Experiments performed on the H100 with FP64 precision. We consider synthetic and real permutations (Section 2.2) in addition to EWA attacks and LP worst-case bounds (Section 2.3).

2.1 Theoretical Description

We show that the complexity of the function describing the decision boundary makes it sensitive to APFPR variability, as shown in the following sections.

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}^L$, $d, L \in \mathbb{N}$, be an arbitrary multiclass classifier where L is the number of classes. Given a datapoint $\mathbf{x} \in \mathbb{R}^d$, the class $1 \dots L$, which the classifier f predicts for \mathbf{x} , is given by $\hat{k}(\mathbf{x}) = \arg \max_i f_i(\mathbf{x})$, where $f_i(\mathbf{x})$ is the i -th component of an array of probabilities called logits. We define a function F at point $\mathbf{x} \in \mathbb{R}^d$ by

$$F(\mathbf{x}) = \max_i f_i(\mathbf{x}) - \max_{i \neq \hat{k}(\mathbf{x})} f_i(\mathbf{x}) \quad (1)$$

F describes the difference between the likelihood of classification for the most probable and the second most probable class. For a given $\mathbf{x} \in \mathbb{R}^d$, the larger the value of $F(\mathbf{x})$, the more confident we are in the prediction given by the classifier. The decision boundary B of a classifier f can then be defined as the set of points f that are equally likely to classify into at least two distinct classes.

$$B = \{\mathbf{x} \in \mathbb{R}^d : F(\mathbf{x}) = 0\} \quad (2)$$

B splits the domain, \mathbb{R}^d into subspaces of similar classification. Given $\mathbf{x} \in \mathbb{R}^d$, and a perturbation $\delta(\mathbf{x}) \in \mathbb{R}^d$ such that $\mathbf{x} + \delta(\mathbf{x}) \in B$, we have that $\mathbf{x} + \delta(\mathbf{x})$ is on the boundary of misclassification. Hence, when considering misclassification, we study properties of the decision boundary B .

Adversarial attacks are small, often imperceptible changes made to input data $\mathbf{x}_{adv} = \mathbf{x} + \delta(\mathbf{x})$, causing the model to misclassify. These perturbations can be viewed as modifications to the input that affect the model’s decision boundary B . Among the most notable adversarial attacks are the gradient based fast gradient signed attack (FGSM) [3] and projected gradient descent (PGD) [11] attack. FGSM generates adversarial examples by adding perturbations in the direction of the gradient of the loss function while PGD can be considered to be an iterative application of FGSM [11]. We also consider a variant of margin-based attacks [23], which we call a *targeted attack*. This targeted attack finds examples closer to the decision boundary B by minimizing the function F . Random attacks adding random noise to the inputs provides a baseline for evaluating model robustness. In general, these attacks are evaluated with an attack scale factor ϵ , such that $\mathbf{x}_{adv} = \mathbf{x} + \epsilon \cdot \delta(\mathbf{x})$ where $\epsilon \geq 0$.

2.2 Synthetic and Real Permutations.

Simple linear classifier: examination of the decision boundary We define a linear classifier by its hyperplane decision boundary:

$$f : \mathbf{x} \in \mathbb{R}^d \longrightarrow \mathbf{1}\{\hat{\mathbf{n}} \cdot \mathbf{x} \geq b\} \quad (3)$$

where $\hat{\mathbf{n}} \in \mathbb{R}^d$ is the normal vector to the hyperplane, $b \in \mathbb{R}$ is the bias, and $\mathbf{1}(\cdot)$ is the indicator function. This classifier assigns inputs to one of two classes based on whether they lie above or below the hyperplane defined by $\hat{\mathbf{n}} \cdot \mathbf{x} = b$. Using the notation from Section 2.1, we express the decision boundary B as:

$$B = \{\mathbf{x} : \hat{\mathbf{n}} \cdot \mathbf{x} - b = 0\} \quad (4)$$

The boundary B is mathematically invariant to any permutation π of the elements $x_i y_i$ in the dot product $\mathbf{x} \cdot \mathbf{n} = \sum_i x_i n_i$. However, run-to-run variability may arise due to APFPR. To simulate the effects of APFPR, we iterate over all possible permutations of the input and normal vector, then compute Eq. 2.2 for points \mathbf{x} on the decision boundary B with $b = 0$. The points \mathbf{x} are sampled from a normal distribution centered at the origin. The resulting distribution of values for $d = 1,000$ and $\hat{\mathbf{n}} = \frac{1}{\sqrt{d(d-1)}} \cdot (d-1, -1, \dots, -1)$ is shown in Fig. 1 (*synthetic permutations curves*). Note we construct the normal vector as such to reduce the space of permutations from $d!$ to d .

The observed distribution exhibits a spread around zero, with a minimum and maximum variation of approximately $\pm 3 \times 10^{-9}$. For the real life case, we perform $N = 1,000$ identical runs for a fixed input. This distribution has a minimum and maximum variation of approximately $\pm 0.9 \times 10^{-9}$. Since the input is fixed and the only source of variation is the accumulation order of floating-point operations, we conclude that APFPR can shift the decision boundary B of the classifier. Furthermore, we show the set of permutations explored in real life identical runs (*run-to-run permutations curves* in Fig. 1) may not cover the set of all possible permutations. Misclassifications may occur as infrequently as once in a thousand identical runs.

To study the effect of input perturbations on classifier robustness, we consider points on the decision boundary $\hat{\mathbf{n}} \cdot \mathbf{x} = b$ and introduce deviations $n \cdot \epsilon + \hat{\mathbf{n}}$, where $\epsilon = 1 \times 10^{-12}$ and $0 \leq n \leq 3000$. As shown in the right panel of Fig. 1, when looping through all possible permutations, classification flips diminish with increasing n , and zero out at $n = 2,400$. We conclude that APFPR cannot affect classification results for inputs far enough from the boundary. We provide a heuristic to identify such "safe" points in Section 2.3. For the real life case, we perform $N = 1,000$ identical runs with similar results, but the percentage of positive predictions zero out much sooner at $n = 700$. The synthetic experiments show that the repeated identical real life runs may not be sufficient to describe robustness, assuming all possible permutations may be explored at runtime. Note the above experiments exhibit similar behavior on FP16, FP32 and FP64 formats, with larger bounds at lower precisions. In Section 2.3, we show how to attack a system to explore permutations that result in misclassification.

2.3 External Workload Attacks and Defense with Learnable Permutations

External Workload Attacks Here we examine the linear decision boundary in Eq. 4 under asynchronous computation on the H100, V100 and Mi250 GPUs. We introduce the EWA, which exploits the impact of background workloads on classification. As studied in detail in Section 4, additional workloads running on the same GPU as inference tasks can affect the ordering of APFPR. We use square matrix multiplications as the background workload, running in a separate process, and determine the optimal matrix size k , that reliably skews classifier outputs.

We use Bayesian optimization with the objective $O(k) = \mathbb{E}[\mathbf{1}(f(x, k), o)]$, where $o \in \{0, 1\}$ is the target output, and $\mathbf{1}(\cdot)$ is the indicator function. Here, we ran the optimization for 100 iterations, with 1,000 experiments per iteration, to find the optimal attack matrix size $k \in \{1,000, 10,000\}$ to flip the classifications into positive classes. Then, we performed 1,000 repeated inferences to test the success of the attack. We highlight that this approach can be generalized to other workloads in addition to matrix multiplication. As shown in the right panel of Fig. 1, all inputs can be reliably skewed toward the desired classification at least 82.7% of the time. The left panel of Fig. 1 shows the positive skewed distribution. The EWA attack is ineffective at $n = 2400$, because no possible configuration of floating point operations results in misclassification (as shown in Section 2.2); this illustrates the effectiveness of the Bayesian approach to find a workload that can exploit any APFPR-based vulnerability. We observe similar behavior on FP16 and FP32 formats.

We find that the EWA optimization convergence behaves similarly across GPU family, however, the optimal matrix size depends on the GPU family (our GitHub repository contains results for the other GPU architectures and datatypes). The trend in inputs and optimal attack matrix size is erratic and we leave an in-depth investigation to future work. Such an investigation would require developing tools to probe the GPU scheduler, which to our knowledge

do not exist [14, 15]. We note that EWA may be inadvertently triggered in cloud systems where GPUs are virtualized and shared. Section 4 further explores this idea by measuring the difference in the scheduling of atomic operations in reductions using black-box testing, both with and without external workloads and analyzing other GPU state factors, including partitioning and power capping.

Learnable Permutation to Find Possible Adversarial Perturbations

We propose a gradient-based optimization technique to find a permutation of floating point operations that causes misclassification. Following Section 2.1, let f be a classifier mapping an input tensor \mathbf{x} to logits, a probability vector of length L , the number of classes. We take the argmax of the logits to obtain the final classification. We require f to include floating point accumulations. Due to APFPR, the output of f depends on a permutation matrix P describing the order reductions, written as $f(P, \mathbf{x})$. In cases where f is composed of multiple functions with the same properties, we parameterize using a set of permutation matrices $\{P_i\}$. The classifier is trained by minimizing a loss function $\mathcal{L}(f, y)$, where y is the ground truth label. To find the $\{P_i\}$ which causes f to misclassify \mathbf{x} , we maximize the prediction error with respect to the permutation perturbation:

$$\begin{aligned} & \text{maximize} && \mathcal{L}(f(\{P_i\}, \mathbf{x}), y) \\ & \text{subject to} && P_i^T P_i = I, \quad i = 1, 2, \dots, L \end{aligned} \tag{5}$$

We use the Gumbel-Softmax technique [5, 12] to create a differentiable approximation of the permutations matrices. By adding Gumbel noise and applying softmax to the matrix, we can use gradient descent to optimize the set $\{P_i\}$ that maximizes the loss function, with the other parameters of f fixed. Next, valid permutation matrices are obtained by solving the linear assignment problem via the Hungarian algorithm [9]. This method, inspired by adversarial attacks (Section 2.1), maximizes error with respect to floating-point operation ordering instead of the input. While the approach does not guarantee misclassification, it provides a more efficient way to find adversarial permutations compared to brute force search.

We now present practical steps to use the LP method: (1) Identify non-deterministic functions by referencing documentation or using a linter like the *torchdet* tool [13], (2) For each non-deterministic function, introduce a permutation matrix P to simulate runtime variations in reductions. For example, consider a fully connected linear layer with a weight matrix w of size $N \times M$ and a bias vector \mathbf{b} of size N , where the intermediate output \mathbf{y} is given by $\mathbf{y} = w^T \mathbf{x} + \mathbf{b}$. To simulate variation in accumulation order, we use the permutation matrix P to permute element-wise products, S_i before reduction where $S_i = \{w_{i0}x_0, \dots, w_{iM}x_M\}$, computing $y_i = \sum_{j=1}^M (P \cdot S_i)_j$. For the linear classifier in Sec 2.2, we introduce the permutation matrix P as follows: $P\hat{\mathbf{n}} \cdot P\mathbf{x}$. (3) Perform gradient descent as specified in Eq. 5, optimizing *only* over permutation matrices. (4) Perform a forward pass and mark any misclassifications to generate a worst-case bound. As shown in the right panel of Fig 1, the LP approach provides a tight bound to the EWA attack. Next, we investigate misclassifications in

GNNs, extending prior work [19] which identified significant run-to-run output variability in GNNs but did not consider misclassification.

3 Non-determinism in Graph Neural Networks

Graph neural networks (GNNs) operate on unordered graph data. For a graph $G = (V, E)$, any permutation of V and E represents the same structure. GNNs learn node and edge representations via *message passing* and *aggregation*, the core operations in most architectures [27]. Since node neighborhoods lack a fixed order, GNNs rely on permutation-invariant aggregation like `add` and `mean` implemented in PyTorch Geometric [1] with `scatter_reduce` functions, which introduce non-determinism due to atomic operations. This, combined with the non-unique representation of graphs, makes GNNs in PyTorch Geometric well-suited for studying APFPR effects.

We investigate run-to-run variability in robustness results, identifying worst-case accuracies with the learnable permutation approach. Additionally, we perform the EWA attack (Section 2.3) to induce misclassifications and evaluate the ability of the LP approach to provide worst-case estimates.

3.1 Experimental Methodology

We study APFPR vulnerability in the GNN architectures: GraphSAGE, GAT, and GCN [4, 8, 24], using the CORA, CiteSeer, and PubMed datasets [16]. These widely used benchmarks evaluate GNN performance in semi-supervised node classification. For each model-dataset pair, we analyze misclassifications due to non-deterministic functions and EWA. We train $N_{\text{train}} = 100$ models for 25 epochs, initializing them identically with fixed randomness from stochastic training and random seed settings. Training models with atomic functions have been shown to produce different weights due to APFPR [19] and we aim to investigate the full training and inference pipeline.

To assess inference variability, we perform $N_{\text{val}} = 10,000$ forward passes on the validation set, with and without atomics. Note the base class of PyTorch GNNs are non-deterministic by default [19]. To provide a deterministic control experiment, we refactor the base class with a deterministic `index_add` operation replacing `scatter_reduce`. However this is *not* a solution to the non-determinism problem in GNNs since significant refactoring is needed to ensure both functional parity and sufficient performance and it is unclear if this is possible. An input is marked misclassified if any iteration produces an incorrect prediction. A large N_{val} is required since misclassifications may only appear after many identical repeated runs, as shown in Section 2.2. We prevent kernel switching, isolating APFPR as the sole source of classification flips. We also predict misclassifications and worst-case accuracy bounds using the LP method (Section 2.3) with 1,000 optimization steps. Graph edges in PyTorch Geometric are permutation-invariant, allowing us to introduce floating-point accumulation sensitivity through permutation matrices in GNN layers when passing the adjacency matrix or `edge_index` variable across layers.

Table 1. Average accuracy (number of correct classifications out of 500) on the CORA dataset for a 10-layer GraphSAGE model under different attacks and attack epsilon values, with standard deviation. "ND" and "D" indicate non-deterministic or deterministic PyTorch settings during inference, respectively. For *ND*, 10,000 inference runs are performed. "LP" refers to a learnable permutation worst-case bound, determined for each input and "EWA" refers to the external workload attack, which succeeds at least 75% of the time on 1,000 repeated runs. We bold experiments which result in misclassifications. All experiments are performed on the H100 with default PyTorch FP32 precision.

Attack	Epsilon	Accuracy D	Accuracy ND	Accuracy LP	Accuracy EWA
None	0	405 ± 9	405 ± 11	402 ± 8	403 ± 8
FGSM	1e − 5	399 ± 8	399 ± 8	397 ± 5	397 ± 6
	1e − 4	397 ± 9	397 ± 9	397 ± 9	397 ± 9
	1e − 3	394 ± 8	394 ± 8	394 ± 8	394 ± 8
	1e − 2	369 ± 9	369 ± 9	369 ± 9	369 ± 9
	1e − 1	340 ± 9	340 ± 10	321 ± 16	328 ± 9
PGD	1e − 5	387 ± 8	387 ± 8	385 ± 9	385 ± 9
	1e − 4	365 ± 9	365 ± 9	365 ± 9	365 ± 9
	1e − 3	348 ± 9	348 ± 9	348 ± 9	348 ± 9
	1e − 2	326 ± 9	326 ± 9	309 ± 8	313 ± 7
	1e − 1	301 ± 9	301 ± 9	287 ± 14	292 ± 15
Random	1e − 5	405 ± 8	405 ± 8	403 ± 10	403 ± 10
	1e − 4	405 ± 9	405 ± 9	405 ± 9	405 ± 9
	1e − 3	405 ± 9	405 ± 9	405 ± 9	405 ± 9
	1e − 2	405 ± 9	405 ± 9	405 ± 9	405 ± 9
	1e − 1	402 ± 9	402 ± 9	383 ± 18	389 ± 20
Targeted	1e − 5	377 ± 8	377 ± 8	375 ± 9	375 ± 9
	1e − 4	365 ± 9	365 ± 9	359 ± 13	361 ± 14
	1e − 3	331 ± 9	331 ± 12	326 ± 5	327 ± 4
	1e − 2	316 ± 9	316 ± 10	298 ± 21	303 ± 21
	1e − 1	293 ± 9	293 ± 9	284 ± 15	288 ± 17

Additionally, we perform the EWA to examine whether external workloads can induce misclassification in a real-world network. As before, we run the Bayesian optimization method for 100 iterations, conducting 1,000 experiments per iteration, to determine the optimal attack matrix size $k \in \{1, 000, 10, 000\}$ for flipping classifications to the second most probable class. We then perform 1,000 repeated inferences to assess attack success, considering it successful if misclassification occurs at least 75% of the time. While this threshold is arbitrary, a reliable attack should consistently induce misclassification.

We validate models and test both unperturbed and adversarial inputs from FGSM, PGD, random, and targeted attacks (Section 2.1).

3.2 Results

The results of experiments for a 10-layer GraphSAGE model on the CORA dataset are shown in Table 1, performed on a H100. Similar behavior was ob-

served across other datasets, models and GPUs (details available on our GitHub). For each adversarial attack method and epsilon value, we report test accuracy as the number of correct classifications (out of 500), with the first row showing results for $\epsilon = 0$ (no attack). Columns labeled *ND* or *D* indicate whether deterministic PyTorch kernels were used during inference, and the LP column represents worst-case upper bounds determined via learned permutation optimization. The EWA column represents the EWA attack, which use naive matrix multiplication workloads. Errors are standard deviations over the 100 trained models.

As expected, accuracies decrease with increasing attack strength ϵ . Toggling PyTorch’s non-deterministic functions has little impact on average accuracy (D and ND columns). However, adversarial accuracy varies significantly at certain epsilon values, indicating that APFPR induces additional misclassifications beyond input perturbations. Notably, large errors occur even at $\epsilon = 0$, showing that non-perturbed inputs are vulnerable to APFPR. EWA reliably misclassifies such inputs, with targeted and PGD attacks being the most affected—leading to adversarial accuracy drops of up to 4.6% (Targeted Attack, $\epsilon = 1e-2$). Random attacks are minimally impacted. The EWA attack works $\geq 75\%$ of the time, which is at least a three order magnitude increase in run-to-run misclassification consistency. EWA achieves convergence in under 80 iterations during the optimization step. Tightening the constraint $\geq 85\%$ makes the EWA ineffective and relaxing it $\leq 20\%$ makes it always effective, producing accuracy values closer to the LP approach. The LP method provides strong worst-case bounds and should be integrated into existing robustness verification tools. A workload as simple and common as a matrix multiplication in the EWA is sufficient to induce misclassifications in GNNs and simpler linear classifiers as in Sec. 2.3. As before, an erratic trend between optimal matrix size and inputs was found. We leave an in depth analysis to future work, requiring the development of novel GPU scheduler probing tools.

4 Impact of GPU State on Order of Reductions

The previous section demonstrates that GPU states can significantly impact ML workloads and should be considered as seriously as more conventional attacks. However, it only provides an indirect measure of the mechanisms behind misclassifications. In this section, we directly measure how the GPU state affects the ordering of asynchronous parallel operations, such as those involving CUDA’s `atomicAdd`. While the exact scheduling of operations on the GPU is not yet fully understood [14, 15], it plays a critical role in determining the order of these operations. To our knowledge, no prior studies have specifically explored how the GPU state influences the ordering of asynchronous operations.

4.1 Methodology

We use the parallel sum algorithm to show how the asynchronous operation order, measured by block index *vs* execution order (BIEO), is influenced by

external loads. The reduction $\sum_i^n x_i$, where $x_i \neq x_j \forall i, j$, $x_i > 0$ double precision floating point numbers (FP64) numbers, runs on a GPU using `atomicAdd`, which has an undefined execution order. To recover the execution order, we track the accumulator updates per block and sort them post-execution. We sum 20 lists of one million uniform FP64 numbers, with and without an additional double-precision matrix-matrix multiplication (DGEMM) workload. Each test runs 10 times to account for system variations. Sorting yields two datasets—reduction only (RO) and with a DGEMM running in a different cuda stream (RDGEMM). The Kendall’s τ correlation [7] measures permutation similarity: K_{RO} for RO and $K_{RO-RDGEMM}$ comparing RO to RDGEMM. Further work probing the GPU scheduler is necessary to investigate reductions of more general arrays.

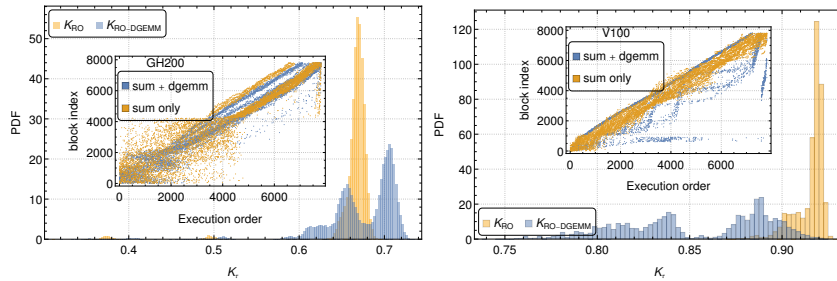


Fig. 2. Left panel: PDF K_{RO} and $K_{RO-RDGEMM}$ on GH200. **Right panel:** $K_{RO-RDGEMM}$ PDF for six MiG configurations on GH200. The inset in the left figure shows BIEO for the RO and RDGEMM workloads with the lowest $K_{RO-RDGEMM}$ correlation.

4.2 Results

This method shows the GPU states effects on the execution order, aiding with verification in non-deterministic settings. We test various GPUs (Section 6), power settings, and partitions, presenting results for GH200 and V100.

Impact of external workloads We calculated K_{RO} and $K_{RO-RDGEMM}$ for the GH200 GPU architecture. As shown in Fig. 2, GH200 exhibits distinct block scheduling and atomic instruction behaviors. The K_{RO} distribution ranges from 0.31 to 0.70, peaking at 0.66, while $K_{RO-RDGEMM}$ ranges from 0.32 to 0.73 with a mean of 0.67. The bimodal distributions for DGEMM workloads differ from the uni-modal distributions for unperturbed reductions, reflecting sensitivity to external workloads. As shown on Fig. 2, the V100 GPU behave differently as K_{RO} and $K_{RO-RDGEMM}$ distributions are narrower than on GH200. The $K_{RO-RDGEMM}$ distribution remains multi-modal but has a more complex structure than on the GH200 GPU. The inset highlights the BIEO snapshot for the pair with the lowest Kendall τ correlation, showing non-sequential block index ordering at first, which later converges into two parallel distributions. These results show that external workloads significantly influence the execution order of `atomicAdd` on GH200, expanding the range of possible ordering permutations.

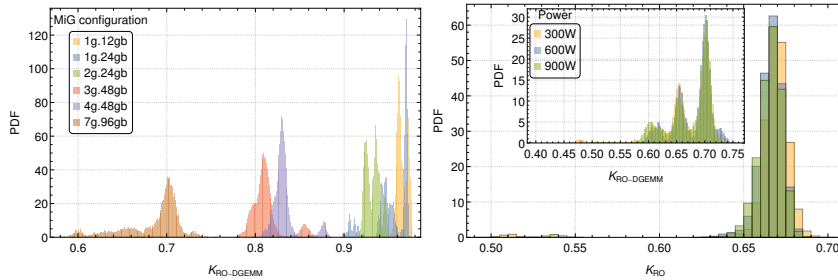


Fig. 3. Left panel: $K_{RO-DGEMM}$ PDF for six different MiG configurations on GH200 showing the effect of resource restrictions on the Kendall τ correlations. **Right panel:** impact of power capping on the Kendall τ correlations.

Impact of MiG Configuration The GH200 supports up to seven MiG partitions, enabling efficient resource sharing. The left panel of Fig. 3 shows $K_{RO-DGEMM}$ across MiG configurations. Higher values occur with fewer SM units and decrease as resources increase while wider distributions in larger GPUs suggest greater permutation variability. Since MiG partitions share the same hardware, each slice’s behavior depends on the runtime configuration. This is especially relevant in virtualized cloud environments.

Impact of Power Research shows that power capping improves energy efficiency and reduces hardware failure rates in HPC and DL workloads [26]. Modern GPUs, like the NVIDIA GH200, include power capping to regulate power draw [26]. While power capping reduces GPU clock speeds, increasing scheduling latency and slowing thread execution, it can also lead to resource contention and impact the atomic operation ordering in compute-intensive tasks. As shown on the right panel of Fig.3 we observed no significant differences in the PDF of K_{RO} and $K_{RO-DGEMM}$, suggesting that power capping does not notably affect instruction ordering.

5 Discussion and Conclusions

We show APFPR has significant impacts on classification accuracy and robustness. We developed a novel black-box Bayesian optimization attack (EWA) to determine properties of additional workloads that reliably result in misclassification (up to 4.6% accuracy decrease, at least 75% of the time). While our current work focused on matrix multiplications as the external background workload, future research will explore more complex workloads in both isolated and cloud environments. Additionally, GPU scheduler probing tools need to be developed to investigate how the EWA and associated workloads impact misclassification.

We introduced the LP approach to efficiently identify permutations that maximize prediction errors, offering significant advantages over brute force search. Our results demonstrate that both run-to-run variability and EWA can be bound by the LP worst-case estimates. Direction for future work involve optimizing LP

further, integrating it into existing robustness verification tools and performing more exhaustive testing over different ML architectures and datasets.

Our examination of GPU system states—including varying workloads, partitions, and power settings—showed significant influence on the ordering of parallel operations across three different GPU models (from NVIDIA and AMD). These findings highlight that testing with a single GPU configuration is insufficient to fully account for non-determinism in model performance. This further emphasizes the value of our LP approach over repeated inferences, which would otherwise require extensive testing across multiple GPU and GPU states.

6 Hardware and Systems Used in Experiments

Tests on GH200 GPUs are run on two separate compute nodes, one running SLE 15 (enterprise) and the other Red Hat Enterprise Linux 9.4 (Plow) with 2 NVIDIA GH200 GPUs and 72-core ARM Neoverse-V2 CPUs. V100 tests use an IBM Power System AC922 (Red Hat 8) with dual Power9 CPUs and six V100 GPUs (16GB HBM2). MI250X tests are performed on an HPE Cray EX (SLE 15) with AMD EPYC CPUs and four MI250X GPUs (eight GCDs per node). H100 tests run on Ubuntu 22.04.06 with two 40GB H100 GPUs and an AMD EPYC 7302 CPU. We use PyTorch 2.4, PyTorch Geometric 2.6 and CUDA 12.0.

References

1. Torch-scatter documentation. <https://pytorch-scatter.readthedocs.io/en/>
2. Gaine, C., Moellic, P.A., Potin, O., Dutertre, J.M.: Fault Injection on Embedded Neural Networks: Impact of a Single Instruction Skip . In: 2023 26th Euromicro Conference on Digital System Design (DSD). p. 317 (2023). <https://doi.org/10.1109/DSD60849.2023.00052>
3. Goodfellow, I.J.: Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572 (2014)
4. Hamilton, W.L., Ying, R., Leskovec, J.: Inductive representation learning on large graphs. In: Proceedings of the 31st International Conference on Neural Information Processing Systems. p. 1025. Curran Associates Inc. (2017)
5. Jang, E., Gu, S., Poole, B.: Categorical reparameterization with gumble-softmax. International Conference on Learning Representations (ICLR) (2017)
6. Jon, P.C., et al: Artificial intelligence for safety-critical systems in industrial and transportation domains: A survey. ACM Computing Surveys **56**, 1 (2023)
7. Kendall, M.G.: A New Measure of Rank Correlation. Biometrika **30**, 81 (1938). <https://doi.org/10.1093/biomet/30.1-2.81>
8. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. In: International Conference on Learning Representations (2017)
9. Kuhn, H.W.: The Hungarian method for the assignment problem. Naval Research Logistics Quarterly **2**, 83 (1955)
10. Lee, Y., et al: Precise extraction of deep learning models via side-channel attacks on edge/endpoint devices. In: Computer Security – ESORICS 2022: 27th European Symposium on Research in Computer Security, Copenhagen, Denmark, September 26–30, 2022, Proceedings, Part III. p. 364 (2022). https://doi.org/10.1007/978-3-031-17143-7_18

11. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. In: International Conference on Learning Representations (2018)
12. Mena, G., Belanger, D., Linderman, S., Snoek, J.: Learning latent permutations with gumbel-sinkhorn networks (2018)
13. MINNERVVA: Torchdet tool. <https://github.com/minnervva/torchdetscan>
14. Olmedo, I.S., Capodieci, N., Martinez, J.L., Marongiu, A., Bertogna, M.: Dissecting the cuda scheduling hierarchy: a performance and predictability perspective. In: 2020 IEEE Real-Time and Embedded Technology and Applications Symposium (RTAS). p. 213 (2020)
15. Otterness, N., Anderson, J.H.: Exploring amd gpu scheduling details by experimenting with “worst practices”. In: Proceedings of the 29th International Conference on Real-Time Networks and Systems. p. 24–34. RTNS '21, Association for Computing Machinery, New York, NY, USA (2021). <https://doi.org/10.1145/3453417.3453432>
16. Prithviraj, S., Galileo Mark, N., Mustafa, B., Lise, G., Brian, G., Tina, E.R.: Collective classification in network data. *AI Magazine* **29**(3), 93–106 (2008)
17. Rabbani, N., Kim, G., Suarez, C., Chen, J.: Applications of machine learning in routine laboratory medicine: Current state and future directions. *Clinical Biochemistry* **103** (2022)
18. Riach, D.: Framework reproducibility: Determinism (d9m). <https://github.com/NVIDIA/framework-reproducibility/blob/master/doc/d9m/README.md>
19. Shanmugavelu, S., Taillefumier, M., Culver, C., Hernandez, O., Coletti, M., Sedova, A.: Impacts of floating-point non-associativity on reproducibility for HPC and deep learning applications (2024). <https://doi.org/10.1109/SCW63240.2024.00028>
20. Summers, C., Dinneen, M.J.: Nondeterminism and instability in neural network optimization. In: International Conference on Machine Learning. pp. 9913–9922. PMLR (2021)
21. Szegedy, C., et al: Intriguing properties of neural networks (2014). <https://doi.org/10.48550/arXiv.1312.6199>
22. Szymanski, N., et al: An autonomous laboratory for the accelerated synthesis of novel materials. *Nature* **624**, 1 (2023)
23. Veerabadran, V., et al: Subtle adversarial image manipulations influence both human and machine perception. *Nature Communications* **14**(1), 4933 (2023). <https://doi.org/10.1038/s41467-023-40499-0>
24. Veličković, P., et al: Graph attention networks. In: International Conference on Learning Representations (2018)
25. Youvan, D.: Parallel precision: The role of gpus in the acceleration of artificial intelligence (2023)
26. Zhao, D., et al: Sustainable supercomputing for AI: GPU power capping at HPC Scale. In: Proceedings of the 2023 ACM Symposium on Cloud Computing. p. 588 (2023)
27. Zhou, J., et al: Graph neural networks: A review of methods and applications. *AI Open* **1**, 57 (2020)