

Meta-Scanner: Detecting Fault Attacks via Scanning FPGA Designs Metadata

Hassan Nassar¹, Jonas Krautter, Lars Bauer¹, Dennis Gnad, Mehdi Tahoori¹, *Fellow, IEEE*,
and Jörg Henkel¹, *Fellow, IEEE*

Abstract—With the rise of the big data, processing in the cloud has become more significant. One method of accelerating applications in the cloud is to use field programmable gate arrays (FPGAs) to provide the needed acceleration for the user-specific applications. Multitenant FPGAs are a solution to increase efficiency. In this case, multiple cloud users upload their accelerator designs to the same FPGA fabric to use them in the cloud. However, multitenant FPGAs are vulnerable to low-level denial-of-service attacks that induce excessive voltage drops using the legitimate configurations. Through such attacks, the availability of the cloud resources to the nonmalicious tenants can be hugely impacted, leading to downtime and thus financial losses to the cloud service provider. In this article, we propose a tool for the offline classification to identify which FPGA designs can be malicious during operation by analysing the metadata of the bitstream generation step. We generate and test 475 FPGA designs that include 38% malicious designs. We identify and extract five relevant features out of the metadata provided from the bitstream generation step. Using ten-fold cross-validation to train a random forest classifier, we achieve an average accuracy of 97.9%. This significantly surpasses the conservative comparison with the state-of-the-art approaches, which stands at 84.0%, as our approach detects stealthy attacks undetectable by the existing methods.

Index Terms—Hardware security, machine learning, reconfigurable logic.

I. INTRODUCTION

FIELD programmable gate arrays (FPGAs) are now heavily utilized as versatile accelerators in the cloud computing domain [1], [2], [3], [4], where the users can realize almost arbitrary circuits on these programmable logic chips. The ever-increasing amount of programmable resources per FPGA chip enables the fine-grained virtualization to optimize efficiency and utilization [5]. Virtualization and multitenancy (multiple

users, i.e., tenants share the resources of the same FPGA) is heavily discussed in [6], [7], and [8]. However, researchers demonstrated unsolved security issues of FPGA multitenancy in the form of remote fault attacks [9], [10], [11], [12]. The attacks have been escalated to the actual cloud devices in the Amazon AWS instances [13], enabling large-scale denial-of-service attacks that can result in financial loss for the cloud service provider (CSP). The attacker causes strong fluctuations in the FPGA's power distribution network (PDN), resulting in its sudden shutdown. The attacker achieves this by implementing several thousands of oscillators on the FPGA [12].

To address these security issues, offline and online countermeasures have been proposed [10], [14], [15], [16], [17], and basic design rule checks are already employed by the industry [13]. Existing offline countermeasures based on the bitstream checking [14], [15], [16] fail to identify the most recent malicious designs. For instance, seemingly benign circuit designs, based on the minor modifications to the AES encryption modules, have been demonstrated as capable of inducing timing faults or causing a denial-of-service [18]. These seemingly benign circuits achieve strong PDN fluctuation through the specific input patterns instead of using simple oscillators. Moreover, it is also possible to build an attack by using multiple malicious tenants in a coordinated way, even though none alone would lead to a successful attack [19]. Thus, attacks are getting more stealthy and are harder to detect.

Two online methods to disable malicious tenant designs during the operation have been proposed [10], [17], but both have restrictions on the type of malicious designs they can prevent and how fast they can do that. Thus, to effectively stop an attack, a hypervisor must know upfront which the tenant could be potentially malicious, as its adversary effects can already become effective a few microseconds after it was deployed. If not, targeting the malicious tenant would take several milliseconds and the attack will be successful [10].

In this work we propose a machine-learning-based classification that can be used offline on the tenant designs before loading them to an FPGA. We show the basic flow of our approach in Fig. 1. In the cloud, a tenant design is compiled into a bitstream with accompanying metadata, such as estimated power consumption. Our classification scans the metadata as input and categorizes the tenant's designs into three categories. These categories correspond to its risk level of becoming a potential threat to the integrity of other tenants or the entire system. Based on the scanner, hypervisors can choose a suitable mapping of tenant designs to different

Manuscript received 12 August 2024; accepted 12 August 2024. This work was supported in part by the “Helmholtz Pilot Program for Core Informatics (kikit)” at Karlsruhe Institute of Technology, and in part by the German Federal Ministry of Education and Research (BMBF) through the Software Campus Project “HE-Trust” under Grant 01IS23066. This article was presented at the International Conference on Hardware/Software Codesign and System Synthesis (CODES + ISSS) 2024 and appeared as part of the ESWEEK-TCAD Special Issue. This article was recommended by Associate Editor S. Dailey. (Corresponding authors: Hassan Nassar; Mehdi Tahoori; Jörg Henkel.)

Hassan Nassar, Dennis Gnad, Mehdi Tahoori, and Jörg Henkel are with the Institute for Computer Engineering, Karlsruhe Institute of Technology, 76131 Karlsruhe, Germany (e-mail: hassan.nassar@kit.edu; mehdi.tahoori@kit.edu; henkel@kit.edu).

Jonas Krautter resides in Wiesbaden, Hessen, Germany.

Lars Bauer resides in Karlsruhe, Baden Württemberg, Germany.

Digital Object Identifier 10.1109/TCAD.2024.3443769

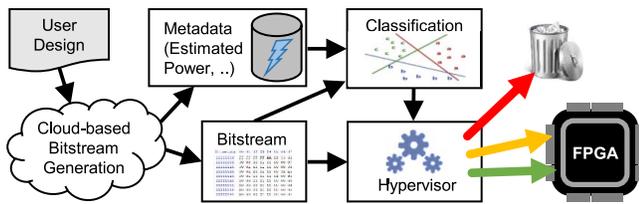


Fig. 1. Basic principle of our proposed meta-scanner and loading flow.

FPGAs that can maximize security through the additional online countermeasures as in [10] and [17].

Tenant designs of the high-risk category (red) are banned from being loaded to any FPGA region. Designs of the low-risk category (green) can be arbitrarily placed together with the other low-risk tenant designs on the same FPGA. The hypervisor would deploy at most a single mid risk (yellow) tenant design on the same FPGA. In this case, existing runtime countermeasures [10], [17] in case of detected malicious activity will be able to disable the yellow tenant design. If we deploy more than one yellow tenant design to the same FPGA, the online tools cannot stop both of them fast enough, in case of an attack.

Altogether, our novel contributions are as follows.

- 1) We propose an offline FPGA design classification in which we identify and extract five relevant features from a tenant design, using the metadata from the bitstream generation step to categorize its risk level.
- 2) Our proposed classifier covers more types of malicious designs than any state-of-the-art solution. It reaches an average cross-validated accuracy of 97.9%, whereas the state-of-the-art checkers only achieve accuracies up to 84.0% in a conservative comparison.
- 3) We generate a comprehensive set of 475 tenant designs based on the malicious and benign logic.¹ We label them using the three risk classes red, yellow, and green.

The organization of the remainder of this article is as follows. We describe the necessary background and state-of-the-art approaches in Section II. Our main contributions, the offline design classification and the required metadata scanning are explained in Section III. Section IV presents the generation of the bitstreams. We present our results and analysis in Section V. In Section VI, we discuss the limitations and advantages of our work. Section VII provides conclusion.

II. BACKGROUND AND RELATED WORK

A major aspect of this work is to focus on the remaining blind spots of the existing countermeasures against the fault attacks in the cloud FPGAs. To get sufficient background, we first explain the current assumptions on the multitenant FPGAs. Then, we detail the existing attacks and their consequences to CSPs. Moreover, we elaborate on which attacks can be performed in the cloud and which further countermeasures are available.

A. Multitenant FPGAs for Cloud Applications

Multitenant FPGAs are a heavily discussed topic. It has interest from both academia [8] as well as industry, e.g., IBM [20]. The idea stems from the fact that the FPGA resources increase and one user, i.e., tenant, might not need to use all the resources on the FPGA. Hence, to increase efficiency, the FPGA can be shared by different tenants. To manage the tenant designs, a static part of the FPGA is used by the CSP to manage the communication and interfaces of the different tenants. The tenant designs reside in a dynamic part, with several accelerator slots that can be used by tenants [5]. The tenant design focuses mainly on the application accelerated by the user [5], [8], [20], [21]. Any memory controllers or PCIe subsystems would not reside in the tenant region. Such components would rather belong to the static design controlled by the CSP to avoid conflicts between the tenants when using the shared resources, such as off-chip RAM. AWS already does this in its commercial single-tenant systems [22].

B. Cloud FPGA Attacks

The interest in multitenant FPGAs sparked security concerns [6], [7], [23]. As in the cloud FPGAs, physical attacks that do not require physical access to the chip become increasingly concerning [12], [24]. In the literature, passive side-channel [24] and active fault attacks [9], [12], [25], [26] are mentioned for the cloud FPGAs. This work solely focuses on the latter.

In such fault attacks, high power-consuming designs cause instability in the PDN. When the attacker uses a large enough design, the whole FPGA or its power supply can crash. This requires manual power cycling to recover the system [12], [25], leading to a major loss of availability. Recently, it was shown that denial-of-service attacks work in commercial FPGA clouds, with only minimal modifications [13]. The authors also show that significant financial loss can be expected for the CSPs, when denial-of-service attacks are performed, leading to longer downtimes of the FPGA infrastructure.

The initial versions [9], [12] of FPGA fault attacks used ring oscillators or other combinational loops for high power consumption. However, these clearly malicious circuits can also be replaced by more stealthy variants, with the first step being synchronous flip-flops [25]. Later, it was shown that the intermittent short-circuits could be caused by certain block RAM access patterns, causing sufficient voltage drop, and even bitflips in configuration memory [27]. Another alternative is “glitch amplification,” which uses a fast-clocked flip-flop with a large output network designed to have many glitches and thus high power consumption [28]. A wider overview of similar circuits optimized for high power consumption is presented in [15].

All of these mentioned circuits use uncommon circuit structures. However, it has also been shown that combining multiple benign synchronous IP modules, e.g., AES, can be used for attacks [18], [29]. Moreover, the attack can be distributed to multiple malicious tenants launching a coordinated attack [19].

¹<https://gitlab.kit.edu/hassan.nassar/Meta-Scanner>

180 C. Offline Countermeasures Against Fault Attacks on 181 Cloud FPGAs

182 From the malicious designs in Section II-B, only some
183 combinational loops are detected by the FPGA CSPs through
184 the typical design-rule checks (DRCs) that are not necessarily
185 meant for security. In the literature, more sophisticated checks
186 have thus been proposed in [14], [15], and [16]. Reverse-
187 engineering is used by [14] and [15] to perform the security
188 checks. They look into detecting patterns to find the malicious
189 elements. Krautter et al. [14] presented a heuristic to check
190 for high-fanout-nets that are often used in attacker designs to
191 toggle large amounts of logic synchronously. Regarding the
192 method presented La et al. [15] used the reverse-engineered
193 bitstream to recreate the netlist. From the netlist, they can
194 find any self-oscillating structures that might escape the DRC
195 done by the CSPs. Moreover, Elnaggar et al. [11] offered a
196 similar approach to ours using ML on the bitstreams. They
197 are limited to work on the full bitstreams, lacking support
198 for the partial bitstreams, and they focus on detecting the
199 self-oscillating structures. They improve over previous works
200 by detecting the hidden malicious designs within the benign
201 designs. Similar to [11] and [16] provides the initial results on
202 training a convolutional neural network (CNN) to detect the
203 self-oscillating structures.

204 However, all these offline countermeasures cannot detect
205 recent malicious designs. As even standard IP-core modules,
206 such as AES and shift registers can be used to provoke
207 crashes [18], [29] because they seem benign. Very recently,
208 Chaudhuri and Chakrabarty [30] and Alrahis et al. [31]
209 showed the initial results for detecting the cryptographical-
210 circuits-based malicious designs. However, they cannot detect
211 many sequential malicious designs, such as shift registers [18]
212 or RAM-based malicious designs [27] which escape detection
213 by all the state-of-the-art solutions. We compare our approach
214 with all the tools from the state-of-the-art in Section V.

215 D. Online Countermeasures Against Fault Attacks on 216 Cloud FPGAs

217 Another mitigation approach is to detect the attacks online,
218 and try to prevent them, i.e., ways for *detection* and for
219 *prevention*. For the detection of attacks, a delay line can be
220 used to detect the voltage drops [12]. By distributing multiple
221 of them, the exact location of the attacker can be found in
222 about 9.9–21.0 μs [17], but some attacks succeed faster than
223 that [10].

224 Preventing attacks can be more challenging, as the FPGAs
225 are not designed to disable an entire region rapidly. When the
226 attack relies on an external clock, a clock disable will be suffi-
227 cient to stop the attack quickly enough [12], [13], [17], but it
228 cannot prevent the attacks with a self-generated clock [10]. To
229 prevent such attacks, *LoopBreaker* can stop attacks at runtime
230 by quickly reconfiguring all the interconnects of the malicious
231 tenant to high impedance in about 1.5 μs [10]. However, due
232 to limitations in the reconfiguration time, *LoopBreaker* needs
233 to know in advance which tenant shall be stopped before the
234 attack even starts. *LoopBreaker* can quickly stop an ongoing
235 attack if and only if that information is available upfront.

III. META-SCANNER: IDENTIFYING MALICIOUS FPGA DESIGNS

236 Our main goal is to develop an offline scanner that allows
237 the CSP to analyse the tenant designs before uploading them
238 to an FPGA. This should be done without a time consuming
239 and extensive netlist analysis, and at the same time, it should
240 be sufficient to complement and assist the existing runtime
241 countermeasures [10], [17]. We classify tenant designs into
242 three categories: 1) high risk (red), 2) mid risk (yellow), and
243 3) low risk (green), which removes the burden from the exist-
244 ing runtime countermeasures to identify the malicious tenant
245 before starting the countermeasure. Our chosen random forest
246 classifier consists of several decision trees. Each decision tree
247 is actually very similar to a simple rule-based inference. The
248 main difference is in finding appropriate thresholds for the
249 rules. The ML part can be seen as an automated way to
250 determine the individual decisions and finding the thresholds
251 during training. This ML training step helps to ensure that
252 the rules are not mistakenly overfitted to the known attacks
253 used for training, but that they remain generic enough to also
254 cover the other attacks. Additionally, it helps to adjusting to
255 novel attack types as soon as they occur, as retraining is an
256 automatic operation.
257
258

A. Threat Model and Assumptions

259 The threat model We target is a cloud scenario with
260 multitenant FPGAs, i.e., multiple tenants share an FPGA
261 in a cloud system with potentially multiple FPGAs. The
262 attacker might rely on *intra-FPGA* coordination, i.e., using
263 multiple regions on a single FPGA *together* to crash the
264 FPGA (see Section II-B). Our focus is mainly on detecting
265 malicious tenant designs. By correctly classifying the risk level
266 of each tenant design, we provide CSPs with the ability to
267 decide whether or not to upload it. We assume that CSPs
268 perform security checks or attestation of the FPGA design
269 through a hypervisor as explained by previous works [32].
270 Moreover, CSPs can combine our risk classification with other
271 data they might have. Usually, CSPs can have access to more
272 information about their users, e.g., their history of previous
273 tenancy on FPGAs. Hence, they may have some trust metric
274 for the users, which is beyond the scope of our work.
275

276 The steps to use our solution are shown in Fig. 1. Normally,
277 a tenant would upload a design as an HDL code or as a netlist
278 to the CSP. The CSP then generates the bitstream and extracts
279 the features (see Section III-C) used by our scanner from the
280 metadata. Then, based on our scanner, the CSP can correctly
281 evaluate the risk category of the bitstream.
282

283 The hypervisor should never upload the red tenants (see
284 Fig. 1), as they will very likely exhibit malicious behavior,
285 whereas the green tenants can always be uploaded, as they
286 are incapable of displaying the malicious behavior. Yellow
287 tenants can be uploaded to an FPGA, but special care must be
288 taken as explained in Section I. When ensuring that at most
289 one yellow tenant is executing on an FPGA, then the online
290 countermeasures like [10] and [17] can aim at the potentially
291 malicious tenant, which allows them to shut it down as soon
292 as it measures any malicious activity. Instead, if two or more
293 yellow tenants were on the same FPGA, it would no longer be

293 known which of them started the malicious activity. Thus, the
 294 online countermeasures would no longer be able to localize
 295 and stop the activity fast enough before a crash occurs.

296 B. Tenant Design Analysis

297 To classify the tenants accordingly, we start by thoroughly
 298 analysing both the malicious- and benign designs (generation
 299 of the dataset of the tenant designs is described in Section IV),
 300 to get an idea of which features would be helpful to detect the
 301 malicious designs. As typical malicious designs aim at trigger-
 302 ing a voltage drop to cause denial-of-service (see Section II-B),
 303 the most straightforward idea is to use the estimated power
 304 consumption of a tenant. However, our analysis shows that this
 305 power estimation is very inaccurate for the earlier published
 306 malicious designs [28] that used highly regular structures
 307 (e.g., mux-based, latch-based, or glitch amplification-based;
 308 see Section II-B). The power estimation alone will not be
 309 enough to classify the malicious designs properly. However,
 310 it is noticeable that the earlier published malicious designs
 311 have a highly regular structure and repetitive elements in their
 312 design, as they are composed of many relatively small building
 313 blocks. We show in Section III-C how to extract and exploit
 314 this property of the bitstream metadata for our classification.

315 Repetitive elements in the bitstream are not always an
 316 indicator of malicious tenants, because simple benign tenant
 317 designs, which have a low power consumption and are mostly
 318 empty. Therefore, they will show a high degree of repetition
 319 in their bitstreams as the unused resources will have similar
 320 configuration data setting them to blank. Hence, these benign
 321 tenants can unintentionally appear like malicious tenants to the
 322 bitstream classifier. The observable repetition is because most
 323 resources in their tenant region are unused. For example, AES
 324 uses only very few DSP blocks. We will have to distinguish the
 325 repetition due to repeated attack blocks from the repetition due
 326 to the repeated unused blocks in the bitstream classification.

327 Complex benign designs like a Bitcoin miner or a cluster of
 328 different big designs have a high estimated power consumption
 329 and a high utilization with a low degree of repetition. It should
 330 be easy to distinguish them from the malicious designs with
 331 highly regular structures. However, malicious designs that are
 332 based on the benign modules (e.g., the AES-based attacks [18]
 333 explained in Section II), also show a high estimated power and
 334 a low degree of repetition, which makes them appear similar
 335 to complex the benign designs.

336 C. Metadata Extraction

337 Our idea is to identify the area utilization of a tenant and
 338 its internal regularity by extracting corresponding properties
 339 directly from its bitstream. Fig. 2 shows the structure of Xilinx
 340 bitstreams. It has headers and trailers for synchronizing the
 341 bitstream upload and the payload. Internally, the main payload
 342 of a bitstream consists of so-called *frames*, i.e., the smallest
 343 reconfigurable unit in an FPGA (in the low kiB range per
 344 frame depending on the FPGA family). For every reconfig-
 345 urable region, the synthesis tools for partially reconfigurable
 346 designs create a so-called *blank* bitstream [shown in Fig. 2(a)]
 347 that reconfigures the region into an empty state.

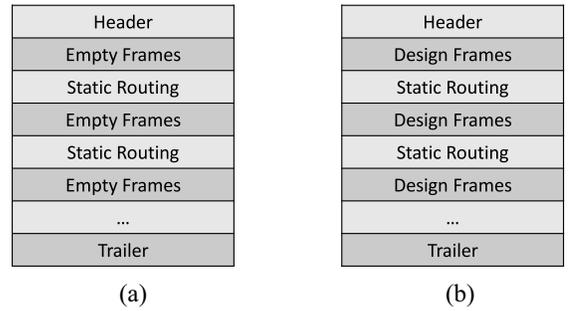


Fig. 2. Bitstream Structure for (a) blank bitstream and (b) design bitstream.

TABLE I
 ANNOTATION OF THE MATHEMATICAL EXPLANATION FOR THE FEATURES

Variable	Explanation
$BitstreamLen$	Number of frames per bitstream
$N_{UFrames}$	Number of unique frames
$N_{BFrames}$	Number of blank (empty) frames
$NonBFrames$	Non Blank Frames

A normal design bitstream for the same region can be seen 348
 in Fig. 2(b). It has the same structure as the blank bitstream. 349
 For unused regions, the frame data is identical to the frame 350
 data of the blank bitstream. Hence, any frame with data 351
 identical to the corresponding frame in the blank bitstream can 352
 be seen as empty. 353

Based on the bitstream structure, we extract five features as 354
 follows. Note that, for the equations, we use the annotation 355
 from Table I. 356

- 1) *Repetition*: The number of nonunique frames. If there 357
 are for instance 100 frames with the identical data, that 358
 adds 100 to the repetition. Nothing is added to the 359
 repetition for an unique frame (i.e., no other frame has 360
 the same data). A higher repetition indicates a higher risk 361
 of self-oscillating structures, as they normally consist of 362
 many repeated frames 363

$$Repetition = BitstreamLen - (N_{UFrames} + N_{BFrames}). \quad 364$$

- 2) *Utilization*: The number of frames different from the 365
 frame data at the same position in the blank bitstream. 366
 This helps to identify the complex designs that use a 367
 large degree of their resources 368

$$Utilization = BitstreamLen - N_{BFrames}. \quad 369$$

- 3) *Average Frame Frequency (AvgFrameFreq)*: We create 370
 a histogram of all the nonblank frames in the bitstream, 371
 i.e., of those frames that are different than the corre- 372
 sponding frame in the blank bitstream. The frequency of 373
 the histogram's bins denotes how many frames belong 374
 to that bin, i.e., how many frames have the same data. 375
 For the AvgFrameFreq, we calculate the average over 376
 the frequencies and divide it by the largest frequency. If 377
 the AvgFrameFreq is near one, it indicates a low degree 378
 of repetition, while if it is close to zero, it indicates a 379
 higher degree of repetition 380

$$AvgFrameFreq = \frac{\text{mean}(\text{hist}(\text{NonBFrames}))}{\text{max}(\text{hist}(\text{NonBFrames}))}. \quad 381$$

4) *Standard Deviation of the Frame Frequency (StdFrameFreq)*: The metric calculates the standard deviation of the frame frequencies and then divides it by the largest frame frequency. This helps identify how much repetition exists. A low deviation means that there is a high degree of repetition and a high deviation means that there is a low degree of repetition

$$\text{StdFrameFreq} = \frac{\text{std}(\text{hist}(\text{NonB}_{\text{Frames}}))}{\max(\text{hist}(\text{NonB}_{\text{Frames}}))}.$$

5) *Estimated Power*: This feature estimates the design’s power consumption. It is the only feature not directly calculated from the bitstream but is reported by the synthesis tools after the design is placed and routed. Note that, for the Amazon cloud, the CSP has access to this information, as the place and route of a tenant design is performed under the control of Amazon (see also Section VI)

$$\text{EstimatedPower} = \text{VivadoPowerEstimation}.$$

Using these five features covers all the important aspects of high utilization, high power, regular structures, and regular structures hidden with some irregularities, which we need for classifying the designs. Overall, they were effective enough to keep our accuracy, recall, and precision around 97%. Initially, we experimented with ten features from the metadata, but through experimentation, we found that the five we use are enough. The five features excluded are the most repeated frame, the number of occurrences of the most repeated frame, the average value of all the bitstream words, and the standard deviation of all the bitstream words. The final five features we use have some overlaps but cover different aspects. *StdFrameFreq* and repetition are somehow correlated. In case of a nonhidden attack, repetition is very powerful to detect the attack while *StdFrameFreq* cannot be of the same strength. However, for the cases where a malicious design is hidden within a benign design, repetition cannot really be used on its own, the *StdFrameFreq* is then more accurate. Therefore, both the features are needed. We evaluate the features’ relevance to our classification problem in Section V-B using the scikit-learn library [33].

D. Proposed Classification

We demonstrate the feasibility of a machine learning approach on the features enlisted in Section III-C by first manually labeling a set of 475 different tenant designs that were tested on a ZCU102 FPGA board (more information in Section IV) according to our three risk classes and evaluating various classifiers on the set. The tenant designs are labeled according to the following principles.

- 1) *RED (High Risk)*: These tenant designs contain actual attack circuits, which we intentionally designed as malicious using different approaches both from [9], [18], [25], [27], and [29]. The hypervisor should never load them to tenant regions on the cloud FPGAs.
- 2) *YELLOW (Mid Risk)*: If a circuit contains a lot of resources and may be used in combination with another similar design on the same FPGA to invoke crashes, we label it as a yellow design. The hypervisor can

permit these designs but requires consideration regarding the mapping into the FPGA regions. Note that, this definition includes completely benign but resource-intensive as well as intentional malicious designs. For instance, additional logic may be added to confuse offline bitstream checker and *hide* the attack, or attackers might use reduced variants of the red designs based on multiple seemingly benign IP modules. Multiple yellow-labeled tenants should not be present at any given time in the FPGA to prevent attacks. If at most a single yellow design is deployed per FPGA, the existing runtime countermeasures [10], [17] will be fast enough to disable it in case of any detected malicious activity (see Section III-A).

- 3) *GREEN (Low Risk)*: Tenant designs from the green category are considered harmless and can be arbitrarily placed into different FPGA regions by the hypervisor. They are neither resource-intensive nor contain known malicious structures, such as self-oscillating circuits. Attacks are highly unlikely even if combined with yellow designs on the same FPGA.

To correctly classify the tenant designs, we use the insights from the bitstream analysis in Section III-B to extract the metadata. By performing the metadata extraction based on the template of the empty tenant region, we can use this metadata to train a lightweight classifier that does not need any complex models to reach far superior results compared to the state-of-the-art as we show in Section V-C.

Based on the recommendations in [34], we evaluate ten-fold cross-validation for different classification methods. We tested a support vector machine (SVM), a multilayer perceptron (MLP), and a random forest classifier. We determined the random forest performed the best on our dataset and used it in all the further experiments. We use the scikit-learn python library [33] to implement the classifier and focus on optimizing the recall for the classification of the red bitstream class by setting the class weights to 200, 30, and 1 for red, yellow, and green, respectively. This approach prevents the misclassification of attack bitstreams into a lower-risk class. Thus, it maximizes the security at the cost of very few lower-risk bitstreams not being loaded to the FPGA.

E. Flow of Using Meta-Scanner

Our proposed Meta-Scanner is easy to use. CSPs will have to deploy a training phase over the existing tenant designs and known malicious designs, then use the trained meta-scanner on any tenant design being uploaded.

1) *Training Phase*: Algorithm 1 summarizes the steps for the training phase by the CSP. It has first to estimate the floorplanning for its different FPGAs to partition them into several tenant regions. The CSP already has several tenant designs from its previous users. For each tenant region, generate the blank bitstream to be used as a reference for the feature extraction. Then, for each tenant design, the bitstream for all the fitting tenant regions has to be generated, and extract the features from the metadata. If no data about whether the design is malicious or not, it has to be uploaded to an FPGA to get the ground truths. Based on the labeled tenant designs, the

Algorithm 1: CSP Classifier Training

Input: List of tenant designs
Output: Trained classifier model

```

ImplementDifferentFloorplans();
foreach tenant region do
    GenerateBlankBitstreams();
    foreach tenant design do
        GenerateBitstream();
        CompareToBlankBitstream();
        ExtractFeatures();
        UploadToFPGAAndGetGroundTruths();
    end
end
FeedLabeledDataToTrainClassifierModel();

```

494 classifier has to be trained to be used to scan the new tenant
495 designs.

496 2) *Scanning Phase:* The trained classifier is continuously
497 used to scan new tenant designs to find out whether they are
498 malicious or not and guide the upload of the tenant designs
499 on the FPGA. The steps of using the scanner are summarized
500 in Algorithm 2. The user usually uploads the tenant design
501 as a synthesized netlist [13]. Therefore, an estimation of
502 the resource needed exists. Based on this estimation, the
503 CSP can choose a suitable tenant region from the floorplan.
504 The CSP generates the bitstream for the tenant region and
505 extracts the features from the metadata. The scanner uses the
506 features extracted to get the label for the design. Based on the
507 label, the tenant design is either banned (red), uploaded with
508 consideration (yellow), or uploaded and trusted (green).

509 IV. DATASET GENERATION

510 To evaluate the effectiveness of our scanner in fulfilling its
511 goal, we generated the dataset of the bitstreams. In Table II
512 we summarize the terminology used to describe the dataset
513 generation.

514 We built a set of bitstreams to extract the metadata and test
515 our solution. The set is based on the 26 basic designs, of which
516 nine are malicious and 17 are benign. We create 475 different
517 tenant designs by configuring, combining, and modifying these
518 26 basic designs. Six of the nine malicious basic designs are
519 from the state-of-the-art mentioned in Section II-B. Moreover,
520 we implement three new malicious designs, similar to the AES
521 malicious design, that we detail later in Section IV-B. The
522 16 benign basic designs are based on the groundhog bench-
523 mark [35], ISCAS benchmark [36], open cores designs [37],
524 Berkeley benchmarks [38], Xilinx HLS tutorials [39], and
525 RISC-V dual core [40]. In addition to these benchmarks,
526 we use some of our developed basic designs, such as JPEG
527 compression/decompression, the secure hash algorithm (SHA),
528 and RSA. We mix and match the basic designs from these
529 benchmarks to build the tenant designs. Table III shows all
530 the basic designs, the benchmarks they originate from, and
531 the frequency of using them in our dataset. Accessing and
532 using the real tenant designs from CSPs is not possible. Even
533 though AWS Marketplace [41] provides FPGA cores, they
534 are typically either simple IP cores meant for integration into

Algorithm 2: Tenant Design Classification and FPGA Deployment

Input: Tenant design netlist
Output: Label

```

// Step 1: Synthesize netlist of the
// tenant design
// Step 2: Estimate a fitting tenant
// region
EstimateTenantRegion();
// Step 3: Perform Place and Route
PerformPlaceAndRoute();
// Step 4: Compare to blank bitstream
CompareToBlankBitstream();
// Step 5: Extract features
ExtractFeatures();
// Step 6: Feed features to classifier
// and get label
// Step 7: Handle label
if Label is RED then
    // Step 7a: Do not upload to FPGA
    Do not upload to FPGA;
end
else if Label is YELLOW then
    // Step 8a: Find suitable FPGA
    FindSuitableFPGA();
    // Step 8b: Upload to FPGA
    UploadToFPGA();
    // Step 8c: Alert online tool
    AlertOnlineTool();
end
else if Label is GREEN then
    // Step 9a: Upload to first fitting
    // FPGA
    Upload to first fitting FPGA;
end
// Step 10: Return label
return Label;

```

TABLE II
TERMINOLOGY USED IN SECTIONS IV AND V

Term	Explanation
Basic Design	HDL code of one module, e.g., DES or JPEG
Tenant Design	One basic design or several of them in a cluster
Tenant region	Area on the FPGA assigned to one tenant
Floorplan	Partitioning the FPGA into different tenant regions
Bitstream	Tenant design in binary, uploaded on the FPGA

larger designs [42], or they are complete systems running
in software that uses hardware IPs. The complete systems
utilize hardware accelerators through an interface without
direct access to the tenant design itself [43]. Therefore, we
rely on the benchmarks as done by [30] and [31] to fulfill
our evaluation, covering a range of applications suitable for
the FPGA acceleration, including neural networks and Bitcoin
mining.

We generate bitstreams for the ZCU102 FPGA board,
utilizing its Xilinx UltraScale+ FPGA for measurements to
establish labeling ground truths. These bitstreams are then

TABLE III
BASIC DESIGNS FOR BITSTREAM GENERATION

Basic Design	Benchmark	#Bitstreams
JPEG	Own Designs	61
RISCV	RISC-V River SoC	15
AVA decoder	Berkeley	42
RSA	Own Designs	46
Cluster of seq. circuits	ISCAS Sequential	64
Serial keyboard	Groundhog	24
PID Controller	Groundhog	45
FIR	Berkeley	28
FFT	Groundhog	40
Bitcoin miner	Opencores	22
AES Attack	Attack from [18]	59
Mux Attack	Attack from [15]	5
Shift register attack	Attack from [18]	20
RAM Attack	Attack from [27]	19
Reed-Solomon attack	Own Designs	19
DES*	Berkeley	25
SHA*	Own Designs	25
Glitch Attack	Attack from [28]	20
Latch Attack	Attack from [15]	20
Neural Network	Opencores	38
Ethernet	Opencores	38
CRC	Opencores	57
SPI	Opencores	57
Manchester encoder	Opencores	38
IIR	Opencores	38
DCT	HLS	25

* Used both maliciously and benignly

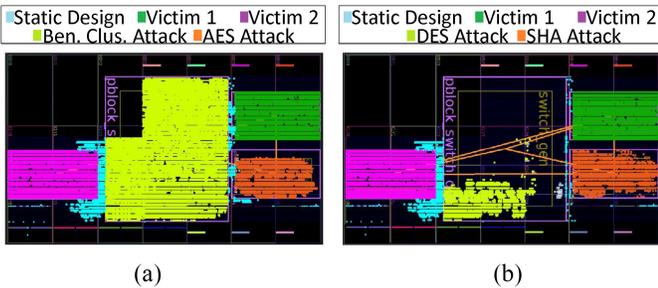


Fig. 3. Floor-planning of tenants. (a) AES and benign cluster coordinated tenant attacks. (b) SHA and DES coordinated tenant attacks.

546 loaded onto the FPGA board. Our focus lies in detecting
547 the success of attacks, which determines the labeling of
548 the bitstreams. The same bitstreams can be used across
549 multiple target FPGA boards, mirroring a cloud scenario from
550 the user's perspective.

551 A. Generating the Tenant Designs

552 We employed various strategies to create the tenant
553 regions. For example, Fig. 3 demonstrates the implementa-
554 tion of the coordinated attacks from the multiple tenants
555 (see Section II-B). The FPGA's floor plan is divided into four
556 regions, with two hosting malicious designs and the other two
557 hosting benign ones. One region utilizes 50% of the resources,
558 while the other three each utilize 15%, leaving 9% for the
559 static design. In the example shown in Fig. 3, the 50% region
560 is positioned in the middle of the FPGA. However, for another
561 floor plan, the 50% region can be placed at the top or bottom of
562 the floor plan, not necessarily in the middle. This contributes
563 to diversifying the bitstreams by avoiding constraining them
564 into fixed regions but instead across several different regions.

A CSP typically maintains several floor plans to accom- 565
modate various types of users. For instance, the 50% tenant 566
region from Fig. 3 can be substituted with two smaller tenant 567
regions, each utilizing 25% of the resources. We employed six 568
different floor plans to generate 24 distinct tenant regions for 569
placing the tenant designs. The sizes of these regions vary, 570
ranging from 50% of the FPGA resources to 15% of the FPGA 571
resources. 572

Not all the tenant designs were used in all the tenant regions 573
as they might not fit into them, i.e., they need more resources 574
than the region provides. Those tenant designs that did not 575
fit were either modified, like changing the RISC-V dual core 576
to a single core, or we diversified the designs further by the 577
following modifications. 578

- 579 1) Mixing them more, e.g., substituting a large FFT 580
module with a smaller PID-controller module and a 581
Manchester encoder. 582
- 583 2) Increasing the repetition within the design, e.g., adding 584
multiple JPEG compression instances after removing a 585
large data encryption standard (DES) module. 586

Moreover, we hide some malicious modules with the benign 585
modules making the attacks stealthier similar to [30]. The 586
generated tenant designs are categorized into 153 green ones, 587
120 red ones, and 177 yellow ones as detailed in Section III-D. 588

B. Implementation of Attacks Based on Benign Constructs 589

We generated malicious tenant designs similar to the AES 590
malicious design from [18] to enrich the dataset. These 591
malicious tenant designs are based on the DES, SHA, and 592
Reed-Solomon as depicted in Fig. 4. The malicious DES- 593
based design in Fig. 4(a) utilizes unrolled DES S-boxes 594
as the fundamental building block. Multiple blocks are 595
interconnected in a chain with adjustable chain lengths to 596
fit the size of the tenant region. The output of each block 597
serves as the input for the subsequent block. The key 598
for each block is computed by XORing the output of the 599
preceding block with the original key. This process amplifies 600
the toggling along the path, thereby increasing the power 601
consumption. 602

The malicious SHA-based design also employs a chain of 603
interconnected SHA subfunctions [shown in Fig. 4(b)]. Each 604
subfunction receives six inputs, which are mixed to produce 605
the various components of the SHA algorithm, resulting in six 606
outputs. The output of one subfunction can be directly con- 607
nected to the next's input, with the chain's length configurable 608
as desired. Note that, only the first input originates from the 609
registers, and no combinational loops are present in the design. 610

As the Reed-Solomon encoder inherently comprises a chain 611
of multiply accumulate operations, the registers between the 612
adder stages are simply removed to transform it into a 613
malicious design [see Fig. 4(c)]. This modification results in 614
a lengthy combinational path, which can be configured as 615
desired. The inputs originate from tenant-internal registers 616
initialized by constants and subsequently inverted in every 617
cycle to enhance toggling. 618

Furthermore, to enhance the difficulty of detection, we 619
explore the concept of hiding these malicious designs among 620

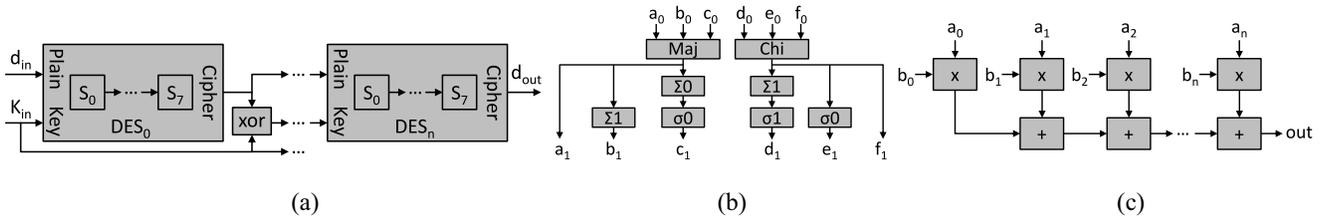


Fig. 4. Implemented attacks, derived from the benign modules. With small modifications, removing sequential elements, and special toggling input patterns, they lead to successful attacks. (a) DES-based attack. (b) SHA-based attack. (c) Reed–Solomon-based attack.

621 the benign ones to evade detection by the current state-of-
 622 the-art solutions. We integrate the malicious designs alongside
 623 a cluster of ISCAS sequential circuits [36]. Consequently,
 624 a bitstream scanner would identify slightly modified benign
 625 designs and encounter additional circuits introducing random-
 626 ness to the structural design. This combined setup presents a
 627 more complicated functionality resembling a standard design,
 628 performing tasks beyond solely cryptographic operations or
 629 encoding.

630 V. EVALUATION

631 We implement the tenant designs using Vivado 2019.1
 632 to evaluate our proposed meta-scanner. The bitstreams were
 633 uploaded to a ZCU102 board. Meta-scanner is implemented
 634 in python and tested on an AMD Ryzen 5 6-Core processor
 635 with 24 GiB main memory.

636 A. Ground Truth of Benign-Based Attacks

637 To label the malicious tenant designs from Section IV-B we
 638 run them on a ZCU102 board to see if they crash the FPGA.
 639 Table IV shows the results. Utilization (%) is based on the total
 640 LUTs available in the ZCU102 FPGA board. Any version of
 641 the malicious designs having the size from Table IV or larger
 642 are labeled as red.

643 Furthermore, we classify smaller malicious designs as
 644 yellow due to their potential for the coordinating attacks,
 645 substantiated by the findings presented in Table IV. Initially,
 646 when both the tenants, SHA and DES are malicious and
 647 deploy weakened versions of their attacks, a coordinated
 648 attack becomes feasible. Second, in scenarios where only
 649 one tenant (AES) is malicious but cannot execute an attack
 650 independently, it can exploit the presence of a resource-
 651 intensive benign tenant. When executed concurrently, the
 652 benign tenant inadvertently facilitates an attack, resulting in a
 653 system crash. Consequently, any benign large design capable
 654 of instigating an attack when combined with the small AES
 655 attack is classified as “yellow.” It should be noted that in
 656 Table IV we show the speed of a crash for the minimum
 657 area. However, using more FPGA resources would cause
 658 faster attacks [10]. Moreover, LoopBreaker [10] can stop an
 659 attack fast only with preselection of the malicious tenant.
 660 Without our tool, LoopBreaker will not be able to identify the
 661 malicious-tenant and would need the lengthy selection step
 662 which requires hundreds of microseconds which is enough for
 663 almost all the attacks to succeed.

664 B. Metadata Features’ Importance

665 As mentioned in Section III-C, we use the scikit-learn
 666 library [33] to evaluate the relevance of our features to the

TABLE IV
 MINIMUM TIME AND UTILIZATION NEEDED FOR ACHIEVING CRASHES

Attack based on	Crash speed	Crash FPGA utilization
AES* [18]	12 μ s	18.5%
Reed-Solomon*	167 μ s	38.7%
DES*	90 μ s	27.0%
SHA*	34 μ s	21.9%
SHA + DES ⁺	60 μ s	14.6% + 18.0%
AES + benign cluster ⁺	>2 Min	13.9% + 34.0%

* attack from single tenant

⁺ attack from multiple coordinated tenants

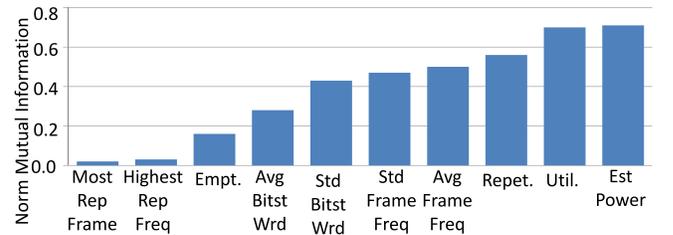


Fig. 5. NMI of each feature individually to the data, 1 is maximum correlation and 0 means no correlation at all.

667 classification problem. In Fig. 5, we show the normalized
 668 mutual information (NMI) between each feature and the data
 669 before classification. NMI is one of the metrics from the
 670 scikit-learn library [33]. It assesses a normalized value with
 671 1 being the highest value (the feature is very relevant to
 672 the classification problem) and 0 being the lowest value (the
 673 feature is not relevant at all to the classification problem).
 674 Utilization and estimated power have the highest NMI values
 675 of around 0.7. The other three features have NMI values
 676 around 0.5. The figure shows that all the metadata features
 677 relate significantly to the classes. Hence, they are all relevant
 678 to the classification problem and can correctly classify the
 679 tenant designs. For the five excluded features mentioned in
 680 Section III-C the NMI score is less than 0.5. Therefore, they
 681 are less suitable for the classification problem and excluding
 682 them is sensible.

683 Attacks rely on the malicious basic designs, e.g., AES-
 684 based, RAM-based, etc. Therefore, we further evaluate their
 685 importance per basic design. To be able to perform this
 686 evaluation, we had to use an one-class classifier method [44].
 687 One-class classifiers belong to the unsupervised learning
 688 approach, where data from only one class is used for training.
 689 The result of the classification is a binary true or false.
 690 For example, if we train an one-class classifier on the DES
 691 malicious designs, it will detect and label them as *true*.
 692 Anything else, even AES or Reed–Solomon malicious designs
 693 would be labeled as *false*.

TABLE V
FEATURE IMPORTANCE PER BASIC DESIGN WHEN USED IN AN ONE CLASSIFIER MODEL

Basic Design	AvgFrameFreq	Repetition	EstimatedPower	StdFrameFreq	Utilization
Artificial Neural Network	2.7%	41.8%	8.9%	18.2%	28.4%
FFT	17.0%	24.0%	15.0%	27.0%	17.0%
RSA	39.4%	11.4%	8.4%	30.6%	10.2%
JPEG	6.9%	32.4%	27.0%	5.6%	28.1%
IIR	13.4%	20.3%	31.1%	15.3%	19.9%
Cluster of s13207/s1494/s9234	14.2%	27.6%	28.0%	7.5%	22.7%
Bitcoin Miner	18.9%	29.4%	16.9%	7.1%	27.7%
Serial keyboard	25.8%	15.8%	16.8%	14.8%	26.8%
SPI	1.0%	45.7%	13.4%	6.3%	33.6%
PID Controller	7.9%	11.8%	41.5%	12.8%	26.0%
DES*	15.4%	11.3%	42.7%	18.9%	11.7%
FIR	28.2%	9.7%	24.3%	18.6%	19.2%
Manchester Encoder	1.0%	36.2%	27.5%	9.8%	25.5%
SHA*	22.0%	6.4%	39.9%	24.6%	7.1%
AES Attack	8.1%	17.1%	44.0%	1.0%	29.8%
AVA decoder	13.5%	19.3%	31.3%	15.6%	20.3%
Ethernet	13.4%	20.3%	31.1%	15.3%	19.9%
RISCV	13.0%	19.8%	31.1%	14.5%	21.6%
Reed Solomon Attack	2.6%	37.3%	28.6%	4.7%	26.9%
CRC	14.1%	21.6%	31.2%	14.1%	19.0%
Latch Attacks	14.7%	14.5%	28.8%	19.9%	22.1%
RAM Attack	12.7%	39.1%	4.5%	29.1%	14.6%
Glitch Attack	19.0%	30.0%	0.0%	32.0%	19.5%
Shift Register Attacks	12.0%	18.6%	12.2%	39.4%	17.8%
Mux Attacks	0.5%	36.2%	9.8%	27.5%	26.0%
DCT	15.3%	26.2%	27.4%	19.2%	11.9%
Average	13.6%	23.6%	25.1%	17.3%	20.1%

* Used both as an attack and as normal module

TABLE VI
RESULTS OF TEN-FOLD CROSS VALIDATION

class	precision	recall	f1score	support	FPR	FNR
GREEN	0.990	0.979	0.984	17.8	0.007	0.020
YELLOW	0.969	0.978	0.972	17.7	0.015	0.016
RED	0.977	0.985	0.979	12.0	0.008	0.021
New RED	1.0	0.963	0.978	1.7	0.000	0.018

The results are presented in Table V. Notably, no feature scored 50% or higher in importance across all the cases. Exceptions were observed where the relevance of features varied among different basic designs. For instance, the Repetition feature held the highest importance for the Reed–Solomon malicious design with a score of 37.3%. Conversely, for the AES malicious design, the power estimation feature scored 44.0%, while utilization scored 25.5%. Additionally, in cases, such as RSA, serial keyboard, and FIR, the average frame frequency feature, typically of low importance exhibited significant relevance. Similarly, estimated power was less critical for the ANN and Bitcoin miner designs, with repetition playing a more substantial role. Furthermore, estimated power played a minor role for several malicious designs like RAM, glitch, mux, and shift registers, reaching 0.0% for the glitch malicious designs. For these malicious designs, the standard deviation of frame frequency gained importance.

C. Performance of the Classifier

The metadata extracted from the bitstream generation is used to train the random forest classifier as described in Section III-D. The training and test data are split randomly by having 10% of the data for testing and the remainder as the

training data. The split is done using the split method from the scikit-learn library [33]. Additionally, we perform ten-fold cross-validation using our 475 bitstreams, and the results are shown in Table VI. The red class has the highest recall and precision to avoid banning legitimate designs and not uploading malicious designs (achieved by fine tuning the class weights as explained in Section III-D). The other two classes (green and yellow) still have high precision and recall and the whole classifier has a mean accuracy of 0.979. Moreover, we ran inference on malicious designs based on our designs from Section IV-B. It had a mean accuracy of 0.95, a precision of 1.0, and a recall of 0.963. For false negatives and positives, Table VI shows that the FPR and FNR are at highest of the value 0.021 which is comparably low. The FPR of the yellow class is roughly the double of the other two classes. This is due to the fact that it is the class in the middle, therefore, a red design will most likely be misclassified as yellow and same for green. For FNR, it is slightly lower for the yellow class than the other two, but in general, it stays low for the all three classes.

Table VII compares our scanner against the five state-of-the-art approaches [11], [14], [15], [16], [30]. As they can only classify into two classes (attack versus no attack), we decided to consider the yellow and green classes as “no attack,” to give them an advantage and to have a conservative comparison. Still, all the state-of-the-art approaches have significantly lower accuracy compared to our scanner. Note that, for the tools from [11] and [30] the tool does not even support partial bitstreams in its current format. However, for a fair comparison, we assume they could be updated to support them. Our scanner is the only tool that detects BRAM short circuit malicious designs and noncryptographic benign-based

TABLE VII
COMPARING OUR SOLUTION TO THE STATE OF THE ART

Metric	Ours	Ref. [14]	Ref. [15]	Ref. [16]	Ref. [11]	Ref. [30]	Ref. [31]
Accuracy	0.979 ± 0.02	0.789	0.756	0.709	0.840	0.825	0.836
Hidden Attacks	✓	✗	✗	✗	✓	✗	✓
Partial Bitstreams	✓	✓	✓	✓	✗	✗	✓
Cryptographic Benign-based Attacks	✓	✗	✗	✗	✗	✓	✓
Non-Cryptographic Benign-based Attacks	✓	✗	✗	✗	✗	✗	✗
Short circuit Attacks	✓	✗	✗	✗	✗	✗	✗
Coordinated Attacks	✓	✗	✗	✗	✗	✗	✗

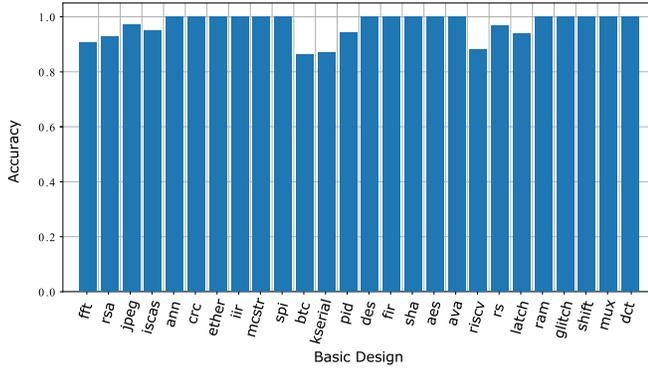


Fig. 6. Mean detection accuracy depending on basic designs.

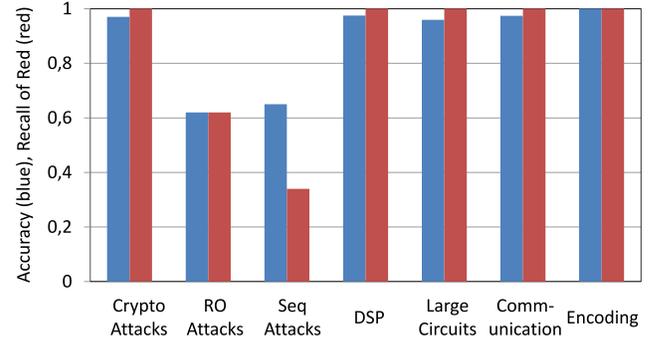


Fig. 7. Accuracy and recall per unseen category.

malicious designs (Reed–Solomon-based and shift-register-based).

Moreover, Fig. 6 shows the accuracy of classifying each basic design to the correct classes. The accuracy is defined as the number of samples correctly classified, divided by the total number of samples used for the inference. Many of the accuracy values are at 1.0, which means that no false positives nor false negatives occur for this basic design. Overall, all the accuracy values are higher than 0.85. DES and SHA (which are both used as the benign designs as well as the malicious designs hidden using the ISCAS circuits) have a high accuracy of 1.0. Hence, our scanner was able to correctly detect hidden malicious designs, and differentiate between using a module for an attack or using it as a true benign design. Moreover, our scanner can detect all the new malicious designs with high accuracy.

Additionally, we evaluate our timing overhead. As mentioned in Section III-A, the CSP performs place and route, feature extraction from the metadata, and scanning (inference of the classifier). Table VIII shows the results of running our scanner on the AMD Ryzen 5 6-Core processor with 24 GiB main memory. On average, place and route for one bitstream needed 27 min, while our feature extraction needs less than 2 s and the inference needs less than 10 ms. Hence, our feature extraction and inference have negligible overhead. The feature extraction takes more time than the inference as it needs to parse the bitstream frame by frame. Moreover, we also measure the time needed for training, our solution needs on average 2 min to train the decision tree.

D. Performance Against Unseen Designs

To complement the classical validation from Section V-C, we use an additional training and test strategy to evaluate the generalization of our classifier. For each basic design b , we perform a training/evaluation experiment, declaring b as

“unseen basic design” and excluding all the bitstreams from the training phase that contains b . This mimics the scenario where a new malicious or benign design emerges that has been used for training. The not-excluded bitstreams are all used to train the model and we test the performance based on the excluded bitstreams. Note that, this evaluation against unseen designs is not performed by any state-of-the-art solution [11], [16], [30], [31]. However, we decided to perform it as an extra step to evaluate the robustness of our scanner.

First, we start by evaluating the case where a full category of designs is unknown. For example, if no cryptography-based attacks were ever used before or if large circuits like neural networks or bitcoin miners are not used before. Fig. 7 shows our results. For most categories, neither recall nor accuracy dropped under 0.9. However, for RO-attacks that use muxes and latches or sequential attacks that use RAM or reed-solomon encoder the accuracy and recall drop. The reason is that these attacks look very similar to the benign small attacks.

We extend our analysis to be even more fine grained. We do it per basic block level Fig. 8 shows (a) the accuracy and (b) recall of the red class for the different unseen basic designs. Recall of the red class is of significant importance, as it shows how well our scanner stops malicious designs from being uploaded. It can be seen from the figures that the scanner’s performance is adequate against the unseen designs with many of them nearly reaching the ideal value of 1 for both the accuracy and recall.

However, there are some outliers. The outliers are analysed and explained in the following, i.e., the Reed–Solomon malicious design, latch malicious design, RAM malicious design, keyboard serial, and the cluster of ISCAS benchmarks. The accuracy and recall are very low for the RAM- and latch-based malicious designs because they are the only malicious designs that do not use any LUTs. Both malicious designs can be implemented using only RAMs or latches, respectively.

900 effort. Our tool works directly on the bitstream and does not
901 require any special maintenance.

902 VII. CONCLUSION

903 In this work, we proposed an meta-scanner, a tool for
904 detecting fault attacks in multitenant cloud FPGA instances.
905 We first analyse the bitstream structure to extract relevant
906 metadata based on them we implemented the classifier for
907 our scanning scheme. By categorizing the client bitstreams
908 into three different risk classes through a machine learn-
909 ing approach, high-risk designs are prevented from being
910 uploaded, whereas the low-risk designs can be mapped to the
911 FPGA regions arbitrarily. Potential attack designs in the mid-
912 risk class can be uploaded, but as long as only a single such
913 design is mapped per FPGA chip, they can be dealt with
914 by existing on-chip countermeasures. Evaluating a random
915 forest classifier on a comprehensive set of 475 different
916 malicious and nonmalicious bitstreams leads to an overall
917 average classification accuracy of 0.979 ± 0.02 , proving the
918 feasibility of our proposed approach. Our solution has a low
919 overhead for training and scanning (inference). Moreover, it
920 can be easily adapted to any new emerging type of attack.

921 REFERENCES

- 922 [1] (Amazon Web Services, Seattle, WA, USA). *EC2 F1 Instances*. 2021.
923 [Online]. Available: <https://aws.amazon.com/ec2/instance-types/f1/>
- 924 [2] (Intel Corp., Santa Clara, CA, USA). *Intel FPGAs Power Acceleration-
925 as-a-Service for Alibaba Cloud*. 2021. [Online]. Available:
926 [https://newsroom.intel.com/news/intel-fpgas-power-acceleration-as-a-
927 service-alibaba-cloud](https://newsroom.intel.com/news/intel-fpgas-power-acceleration-as-a-service-alibaba-cloud)
- 928 [3] (Huawei Cloud, Guangzhou, China). *FPGA Accelerated Cloud
929 Server*. 2021. [Online]. Available: [https://www.huaweicloud.com/en-
930 us/product/fcs.html](https://www.huaweicloud.com/en-us/product/fcs.html)
- 931 [4] (Open Telekom Cloud, Frankfurt, Germany). *Faster Computing:
932 FPGA Hardware Acceleration for the Open Telekom Cloud*. 2021.
933 [Online]. Available: [https://open-telekom-cloud.com/en/blog/product-
934 news/fpga-closed-beta](https://open-telekom-cloud.com/en/blog/product-news/fpga-closed-beta)
- 935 [5] S. A. Fahmy, K. Vipin, and S. Shreejith, "Virtualized FPGA accelerators
936 for efficient cloud computing," in *Proc. IEEE 7th Int. Conf. Cloud
937 Comput. Technol. Sci.*, 2015, pp. 430–435.
- 938 [6] Y. Luo and X. Xu, "HILL: A hardware isolation framework against
939 information leakage on multi-tenant FPGA long-wires," in *Proc. ICFPT*,
940 2019, pp. 331–334.
- 941 [7] G. Dessouky, A.-R. Sadeghi, and S. Zeitouni, "SoK: Secure FPGA multi-
942 tenancy in the cloud: Challenges and opportunities," in *Proc. IEEE Eur.
943 Symp. Security Privacy*, 2021, pp. 487–506.
- 944 [8] M. G. Jordan, G. Korol, M. B. Rutzig, and A. C. S. Beck, "MUTEKO:
945 A framework for collaborative allocation in CPU-FPGA multi-tenant
946 environments," in *Proc. SBCCI*, 2021, pp. 1–6.
- 947 [9] J. Krautter, D. R. E. Gnad, and M. B. Tahoori, "FPGAhammer: Remote
948 voltage fault attacks on shared FPGAs, suitable for DFA on AES," *IACR
949 Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2018, no. 3, pp. 44–68, 2018.
- 950 [10] H. Nassar, H. AlZughbi, D. Gnad, L. Bauer, M. Tahoori, and J. Henkel,
951 "LoopBreaker: Disabling interconnects to mitigate voltage-based attacks
952 in multi-tenant FPGAs," in *Proc. ICCAD*, 2021, pp. 1–9.
- 953 [11] R. Elnaggar, J. Chaudhuri, R. Karri, and K. Chakrabarty, "Learning
954 malicious circuits in FPGA bitstreams," *IEEE Trans. Comput.-Aided
955 Design Integr. Circuits Syst.*, vol. 42, no. 3, pp. 726–739, Mar. 2023.
- 956 [12] D. R. E. Gnad, F. Oboril, and M. B. Tahoori, "Voltage drop-based fault
957 attacks on FPGAs using valid bitstreams," in *Proc. FPL*, 2017, pp. 1–7.
- 958 [13] T. La, K. Pham, J. Powell, and D. Koch, "Denial-of-service on FPGA-
959 based cloud infrastructures—Attack and defense," *IACR Trans. Cryptogr.
960 Hardw. Embed. Syst.*, vol. 2021, no. 3, pp. 441–464, 2021.
- 961 [14] J. Krautter, D. R. E. Gnad, and M. B. Tahoori, "Mitigating electrical-
962 level attacks towards secure multi-tenant FPGAs in the cloud," *ACM
963 Trans. Reconfig. Technol. Syst.*, vol. 12, no. 3, pp. 1–26, 2019.
- 964 [15] T. M. La, K. Matas, N. Grunchevski, K. D. Pham, and D. Koch,
965 "FPGADefender: Malicious self-oscillator scanning for Xilinx
966 UltraScale+ FPGAs," *ACM Trans. Reconfig. Technol. Syst.*, vol. 13,
967 no. 3, pp. 1–31, 2020.
- [16] J. Chaudhuri and K. Chakrabarty, "Detection of malicious FPGA
968 bitstreams using CNN-based learning*," in *Proc. ETS*, 2022, pp. 1–2.
969
- [17] G. Provelengios, D. Holcomb, and R. Tessier, "Mitigating voltage attacks
970 in multi-tenant FPGAs," *ACM Trans. Reconfig. Technol. Syst.*, vol. 14,
971 no. 2, pp. 1–24, 2021.
972
- [18] G. Provelengios, D. Holcomb, and R. Tessier, "Power wasting circuits
973 for cloud FPGA attacks," in *Proc. FPL*, 2020, pp. 231–235.
974
- [19] H. Nassar, P. Machauer, D. R. E. Gnad, L. Bauer, M. B. Tahoori, and
975 J. Henkel, "Covert-hammer: Coordinating power-hammering on multi-
976 tenant FPGAs via covert channels," in *Proc. ACM/SIGDA ISFPGA*,
977 2024, p. 43.
978
- [20] Z. Zhu, A. X. Liu, F. Zhang, and F. Chen, "FPGA resource pooling
979 in cloud computing," *IEEE Trans. Cloud Comput.*, vol. 9, no. 2,
980 pp. 610–626, Apr.–Jun. 2021.
981
- [21] M. Paolino, S. Pinnetter, and D. Raho, "FPGA virtualization with accel-
982 erators overcommitment for network function virtualization," in *Proc.
983 ReConFig*, 2017., pp. 1–6.
984
- [22] *Introduction of F1 Development Environment*, Amazon Web Services,
985 Seattle, WA, USA, 2021.
986
- [23] H. Englund and N. Lindskog, "Secure acceleration on cloud-based
987 FPGAs—FPGA enclaves," in *Proc. IEEE IPDPSW*, 2020, pp. 119–122.
988
- [24] F. Schellenberg, D. R. Gnad, A. Moradi, and M. B. Tahoori, "An inside
989 job: Remote power analysis attacks on FPGAs," in *Proc. DATE*, 2018,
990 pp. 1111–1116.
991
- [25] T. Sugawara, K. Sakiyama, S. Nashimoto, D. Suzuki, and T. Nagatsuka,
992 "Oscillator without a combinatorial loop and its threat to FPGA in data
993 centre," *IET Electron. Lett.*, vol. 55, no. 11, pp. 640–642, 2019.
994
- [26] A. Boutros, M. Hall, N. Papernot, and V. Betz, "Neighbors from
995 hell: Voltage attacks against deep learning accelerators on multi-tenant
996 FPGAs," in *Proc. ICFPT*, 2020., pp. 103–111.
997
- [27] M. M. Alam, S. Tajik, F. Ganji, M. Tehranipoor, and D. Forte, "RAM-
998 jam: Remote temperature and voltage fault attack on FPGAs using
999 memory collisions," in *Proc. FDTIC*, 2019, pp. 48–55.
1000
- [28] K. Matas, T. M. La, K. D. Pham, and D. Koch, "Power-hammering
1001 through glitch amplification—Attacks and mitigation," in *Proc. FCCM*,
1002 2020, pp. 65–69.
1003
- [29] J. Krautter, D. R. E. Gnad, and M. B. Tahoori, "Remote and
1004 stealthy fault attacks on virtualized FPGAs," in *Proc. DATE*, 2021,
1005 pp. 1632–1637.
1006
- [30] J. Chaudhuri and K. Chakrabarty, "Diagnosis of malicious bitstreams in
1007 cloud computing FPGAs," *IEEE Trans. Comput.-Aided Design Integr.
1008 Circuits Syst.*, vol. 42, no. 11, pp. 3651–3664, Nov. 2023.
1009
- [31] L. Alrahis et al., "MaliGNNoma: GNN-based malicious circuit
1010 classifier for secure cloud FPGAs," in *Proc. IEEE HOST*, 2024,
1011 pp. 383–393.
1012
- [32] S. Zeitouni, J. Vliegen, T. Frassetto, D. Koch, A.-R. Sadeghi, and
1013 N. Mentens, "Trusted configuration in cloud FPGAs," in *Proc. FCCM*,
1014 2021, pp. 233–241.
1015
- [33] F. Pedregosa et al., "Scikit-learn: Machine learning in python," *J. Mach.
1016 Learn. Res.*, vol. 12, no. 85, pp. 2825–2830, 2011.
1017
- [34] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical
1018 Learning*. New York, NY, USA: Springer, 2009.
1019
- [35] P. Jamieson, T. Becker, P. Y. K. Cheung, W. Luk, T. Rissa, and
1020 T. Pitkänen, "Benchmarking and evaluating reconfigurable architectures
1021 targeting the mobile domain," *ACM Trans. Design Autom. Electron.
1022 Syst.*, vol. 15, no. 12, pp. 1–24, 2010.
1023
- [36] F. Brglez, D. Bryan, and K. Kozminski, "Combinational profiles of
1024 sequential benchmark circuits," in *Proc. ISCAS*, 1989, pp. 1929–1934.
1025
- [37] (OpenCores, Amsterdam, The Netherlands). *OpenCores the Reference
1026 Community for Free and Open Source Gateway IP Cores*. 1999.
1027 [Online]. Available: <https://opencores.org>
1028
- [38] J. Pistorius, M. Hutton, A. Mishchenko, and R. Brayton, "Benchmarking
1029 method and designs targeting logic synthesis for FPGAs," in *Proc. IWLS*,
1030 2007, pp. 1–8.
1031
- [39] *Vitis HLS Tutorial (v2023.1)*, Xilinx, Inc., San Jose, CA, USA, 2023.
1032
- [40] A. Waterman, Y. Lee, D. A. Patterson, and K. Asanovic, "The RISC-V
1033 instruction set manual," Dept. Electr. Eng. Comput. Sci., Univ. California
1034 Berkeley, Berkeley, CA, USA, UCB/EECS-2011-62, Rep. 2014.
1035
- [41] "AWS Marketplace." 2022. [Online]. Available:
1036 <https://aws.amazon.com/marketplace>
1037
- [42] "VFI—Accelerator FPGA." 2022. [Online]. Available: [https://aws.
1038 amazon.com/marketplace/pp/prodview-fboudtmufjkum](https://aws.amazon.com/marketplace/pp/prodview-fboudtmufjkum)
1039
- [43] *Hugenomic Nanopolish*, Huxelerate, Milan, Italy, 2022.
1040
- [44] S. S. Khan and M. G. Madden, "One-class classification: Taxonomy
1041 of study and review of techniques," *Knowl. Eng. Rev.*, vol. 29, no. 3,
1042 pp. 345–374, 2014.
1043