

# Verification of Cyber Emulation Experiments Through Virtual Machine and Host Metrics

Jamie Thorpe and Laura P. Swiler and Seth Hanson and Gerardo Cruz and Thomas Tarman  
{jthorpe,lpswile,shanson,gcruz,tdtarma}@sandia.gov  
Sandia National Laboratories  
Albuquerque, New Mexico, USA

Trevor Rollins and Bert Debusschere  
tkrollins6@gmail.com,bjdebus@sandia.gov  
Sandia National Laboratories  
Livermore, California, USA

## ABSTRACT

Virtual machine emulation environments provide ideal testbeds for cybersecurity evaluations because they run real software binaries in a scalable, offline test setting that is suitable for assessing the impacts of software security flaws on the system. Verification of such emulations determines whether the environment is working as intended. Verification can focus on various aspects such as timing realism, traffic realism, and resource realism. In this paper, we study resource realism and issues associated with virtual machine resource utilization. We examine telemetry metrics gathered from a series of structured experiments which involve large numbers of parallel emulations meant to oversubscribe resources at some point. We present an approach to use telemetry metrics for emulation verification, and we demonstrate this approach on two cyber scenarios. Descriptions of the experimental configurations are provided along with a detailed discussion of statistical tests used to compare telemetry metrics. Results demonstrate the potential for a structured experimental framework, combined with statistical analysis of telemetry metrics, to support emulation verification. We conclude with comments on generalizability and potential future work.

## CCS CONCEPTS

• Security and privacy → Formal methods and theory of security.

## KEYWORDS

cyber experimentation, system emulation, model verification

## ACM Reference Format:

Jamie Thorpe and Laura P. Swiler and Seth Hanson and Gerardo Cruz and Thomas Tarman and Trevor Rollins and Bert Debusschere. 2022. Verification of Cyber Emulation Experiments Through Virtual Machine and Host Metrics. In *Proceedings of Cyber Security Experimentation and Test Workshop (CSET'22)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/XXXXX.123456>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

CSET'22, August 2022, Virtual Conference

© 2022 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/XXXXX.123456>

## 1 INTRODUCTION

Cyber networks are commonly modeled using network emulation [12], which runs actual software binaries (e.g. operating systems and applications) on virtualized hardware, offline and separated from real operational systems. A number of cyber emulation testbeds have been developed to provide a platform for test and evaluation, cybersecurity investigation, research, and training. [17, 27]. Such testbeds include LARIAT [26], Emulab [28], DETER [7, 24], and DARPA's National Cyber Range [13]. Experimental frameworks such as DEW [23] and SCORCH [14] have been developed to run structured experimental scenarios on these testbeds.

When using emulation, it is important to determine whether the emulation environment is working as intended, also called *verification* [1] [25]. Verification of systems in general typically involves software testing and quality assurance. A unique aspect of cyber emulation in particular is assessing the performance of the experiment in the virtualized environment and determining whether there are sufficient resources available to run the experiment properly. If there are not, the virtualized components may produce experimental artifacts that cause the outcomes to be incorrect or unrepresentative of the system being modeled.

The goal of this study is two-fold: (1) develop and exercise a process for verifying cyber emulation environments, and (2) identify metrics that can indicate when there are insufficient resources to reliably run an emulation. In this work, we refer to these metrics as *telemetry metrics*, following the usage of this phrase from Google [5], Microsoft [3], Intel [4] and others [6]. These companies use telemetry in the context of network monitoring metrics [5], virtual machine resource usage [3], and application monitoring. In this paper, we focus on telemetry metrics relating to the performance of virtual machines which are used in a cyber emulation study and the physical machine hosting that study, specifically for an emulation tool called *minimega* [22].

This study was performed on two different scenarios, described later in this paper. These scenarios are specific examples meant to demonstrate our proposed general approach to resource verification. While there has been some related research in emulation verification and validation, there is generally no common framework, set of standards, or evaluation metrics to use when setting up an emulation; it is left to the analyst running the emulation model to determine whether the emulation ran correctly. We hope to make this process more consistent and objective by providing a process for resource verification and a method for identifying metrics which can be used to detect anomalous experimental results.

The organization of this paper is as follows: Section 2 provides an overview of related research. Section 3 summarizes our approach to

investigating telemetry metrics for emulation verification. Section 4 outlines the experimental configurations used in both scenarios of this study. Section 5 presents a discussion of statistical metrics used to compare telemetry metrics. Section 6 presents the results for both scenarios. Section 7 discusses future work.

## 2 RELATED WORK

Verification and validation (V&V) of emulation environments often goes hand in hand with V&V for software that can be run in those environments. Simulation or emulation models can be used for initial software V&V testing. However, as discussed by Zheng et al. [32] in the context of cyber-physical systems, successful V&V in simulation/emulation does not always translate to real-world systems. When this happens, we can say that the software under test failed verification and/or validation. However, we may also be able to say that the modeling environment itself may not have been successfully verified or validated. Underlying issues in the fidelity of the construction of the system model or in the manner in which it was run could cause results from the model to not be fully predictive of real-world performance.

Before turning our attention to the focus of this paper, verification, it is important to address *validation*. Validation addresses whether the behavior of a lower-fidelity model is “close enough” to the behavior of a high-fidelity model to be an acceptable and adequate representation. Brown et al. [8] investigated network validity in simulation models compared to emulation models, with the goal of determining whether simulations could be an effective and more time- and cost-efficient method for cyber experimentation in certain use cases. Crussell et al. also investigate validation of emulation models to physical systems [11], identifying behavior at three levels of abstraction: application, operating system, and network. A unique aspect of their work was the use of Markov models to compare patterns of system call orderings. While we acknowledge validation is also an important area of research, the focus of our study is resource verification.

It is useful to consider resource verification as related to resource allocation, where the goal is to ensure that processes (e.g., virtualized networks) have sufficient resources available to them. These challenges are often explored in the context of testing allocation algorithms for cloud environments [9], [19]. However, resource verification in the way that we define it is not a goal of that research. Similarly, work by Hibler et al. [16] considered resource allocation in the development of the emulation testbeds Emulab. Emulab collects resource metrics (e.g., CPU, memory, bandwidth) from interactive experiments over time, enabling resources to be reallocated adaptively to better support the virtual network. The authors noted differences in time to perform tasks as more resources were used by more virtual nodes hosted on the same physical node. Resource verification would next ask whether these noted changes in resource use had significant impact on the key outcomes of the experiment.

In his Ph.D. thesis [15], Brandon Heller developed the idea of “network invariants.” These are properties that can be used to verify timing errors in the emulation. The network invariants should be universal and apply regardless of the experiment being run, the topology of the experiment, the platform upon which the emulation

is run, or the emulation code being used. In this paper, we expand upon Heller’s idea of network invariants and consider experimental scenarios where results depend less on timing errors and more on available resources.

It is also important to consider whether periodic polling to track the metrics, a form of “health monitoring”, affects the behavior of the emulation itself by using resources. Jia et al. [18] referred to this issue as a weak form of the Heisenberg uncertainty principle for measuring host resource usage. They studied virtualized environments where hundreds of VMs were running on one physical host. They concluded that frequent polling (sampling intervals between 0.001 and 0.1 seconds) did affect the performance of the emulated environment. In the work we present in this paper, the polling was much slower, at 5 or 10 second intervals, and our experimentation suggests this frequency of polling did not introduce artifacts into the results.

## 3 SUMMARY OF APPROACH

The first goal of this study is to develop and exercise a process for resource verification. The process we developed is as follows:

- (1) Plan a series of experiments to stress the physical resources of your system. We refer to a series of these experiments as “runs”. Runs are defined by different emulation configurations which perform increasing numbers of experiments in parallel. Each run will have the same total number of “replicates”, or completed instances of the emulation experiment.
- (2) Configure the collection of telemetry data from the replicates. In addition to the experiment results (“quantities of interest”) for each replicate, telemetry metrics are collected to help answer questions about the emulation itself, and to verify that the emulation had sufficient resources to run.
- (3) Perform the planned experiments, increasingly stressing the physical host to the point that emulation artifacts are generated. Purposefully pushing the physical resources to oversubscription allowed us to address the secondary goal of the study: to develop indicators that could help detect unreliable experiments run in oversubscribed conditions.
- (4) Observe changes to the emulated scenario’s quantity of interest as resource conditions vary. If the quantity of interest changes drastically as more and more resources are utilized on the physical host, then we know that there are likely emulation artifacts affecting the results.

In addition to establishing this verification process, the second goal of this study is to identify telemetry metrics that can indicate when the results of an experiment are likely unreliable. We define an “indicator” as a combination of a telemetry metric and a threshold which that metric should not cross. We wish to identify indicators where violation of the threshold during an experiment is highly correlated with unreliable results in the quantity of interest. We hypothesize that this indicator can then be used to detect experiments with unreliable results even when the level of resource subscription is not known.

We evaluate our approach using two scenarios, described further in Sections 4.2 and 4.3.

## 4 EXPERIMENTAL SETUP

### 4.1 Emulation Infrastructure

- **minimega and SCORCH**

The emulation platform we used to run the emulations is minimega [22]. We also leveraged SCORCH, a framework that handles scenario orchestration for minimega. SCORCH allows specific emulation experiments to be defined in a highly modular way. [14]. Modules called “components” define a variety of aspects of the experiment, including the network topology, the scenario being run, and tools for data handling. Components can also be configured to feed resulting data directly to a database for later analysis.

- **Over-Committing Resources**

Straining the available physical resources will not only help test our verification process, but also give us sufficient data to develop telemetry-based indicators of unreliable experiments. We oversubscribe resources by forcing the physical host to perform increasingly more work in parallel. This is implemented using SCORCH’s “namespace” parameter. Namespaces can be thought of as isolated copies of the experiment environment, as specified by the scenario being run. Multiple namespaces can be started in parallel, each with its own copy of the experiment environment, and each run its own series of experiments without affecting the actions of the other namespaces. By increasing the number of namespaces for subsequent runs, we hope to reach a point of resource over-subscription.

- **Data Collection**

In order to establish indicators of unreliable experiments, data must be collected from the emulations. Existing SCORCH components which define the basic experiment may be modified to also collect data or to collect data in new ways. Alternately, new SCORCH components may be written to collect specific data from new sources. Components collect data at a given rate and output results which can be parsed and stored locally or in a database. For this study, collected data included telemetry from the emulated VMs and the physical host running the emulations, as well as the results of the experiment scenario.

### 4.2 Scanning and Detection Scenario

The scanning and detection scenario is extensively documented in [30] and [29]. It consists of an attacker conducting reconnaissance using a scanning tool called nmap and is detected by the intrusion detection system called snort [10].

Two modes of stochasticity, port scan order and rate of packet dropping, are removed from this scenario for the purpose of this study – the scanning order is fixed and the drop rate is set to zero. The quantity of interest for this scenario is the time to detection, or “alert” time.

### 4.3 Command and Control Scenario

The Command and Control scenario [31] emulates a network under attack by malware with a command and control (C2) communication pattern, such as Emotet [2]. C2 involves a malicious actor

leveraging a covert communication channel to communicate with infected hosts within a network. A snort intrusion detection system on this network attempts to detect the malicious C2.

The Command and Control scenario was emulated using a simplified network topology. This reduced emulation complexity while still allowing us to track the number of alerts over time. The emulated network consisted of a VM to generate background traffic, a VM to generate malware traffic, a single traffic server VM, and a single VM to host snort. There were several tunable parameters for this scenario, but all parameters were set to constant values for the duration of the verification study. There were 1000 benign packets per second, and 1 in 1000 packets contained the malware signature that snort would identify. There were 20 malicious packets per second, and 3 in 20 (150 in 1000) contained the malware signature.

The quantities of interest for this scenario are the number of alerts received by time  $t = 1, 5, 10,$  and 16 seconds. Because this scenario has multiple quantities of interest, the statistical analysis differed from that of the Scanning and Detection scenario. See Section 5.2 for details.

## 5 METHODS

### 5.1 Statistical Comparisons of Experimental Results

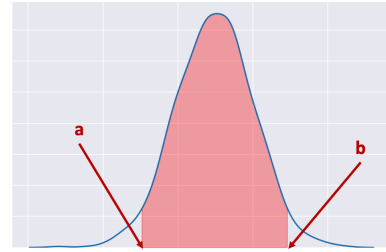
As discussed in the scenario descriptions above, each emulation scenario has a quantity of interest that is being monitored (e.g. time to alert in the scanning and detection scenario). The telemetry metrics such as CPU utilization and system load are also being monitored as a function of time (e.g. recorded every second). The distribution of the quantity of interest is used to compare similarity of experimental outcomes at each namespace case. The assumption is made that the 1-namespace case should be the expected results of these experiments (as these are least likely to oversubscribe the physical resources), and the distribution of the 1-namespace quantities is taken as the ground truth. Two statistical comparisons are the Tukey multiple comparison test and the ratio acceptance test as described below. Additional statistical tests will be the subject of future work.

*5.1.1 Tukey Multiple Comparison Test.* Tukey’s multiple comparison test is used to determine which means amongst a set of means differ from the rest. [21] A mean value in this context is the average of a number of experimental replicate results. For example, the Tukey test allows us to compare the alert time mean at each namespace case to all other namespace cases. That is, for a set of  $r$  means, the Tukey test considers all possible pairwise differences:  $\mu_i - \mu_j$  for  $i \neq j$  and  $i, j = 1 \dots r$ . The Tukey test produces a p-value for each comparison, with a high p-value indicating that the null hypothesis:  $\mu_1 = \mu_2$  cannot be rejected. Thus, a high p-value indicates that our distributions of alert times are similar, indicating that the experiment produced comparable results in the two cases.

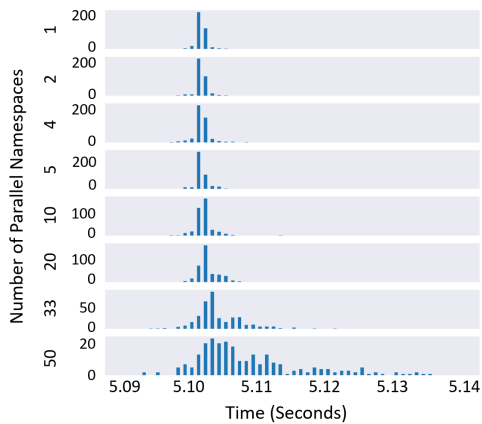
An example of how the alert times change in the scanning and detection scenario is shown in Figure 1. The set of histograms in subplot (a) show the systematic change in the alert time distribution as the number of namespaces increases from one to 50. The 50 namespace experiment has a wider, longer tailed distribution and the mean is also larger. The matrix of Tukey test values in subplot

(b) of Figure 1 shows that the alert times in the tests run on 20, 33, or 50 namespaces would be considered significantly different from the one namespace results, as evidenced by the small p-values highlighted by the red colors. Thus, the Tukey test serves to quickly summarize differences in mean alert times across varying namespaces. Note that p-values appearing as “0.00” in the Tukey figures throughout this paper indicate a non-zero p-value which is less than 0.005.

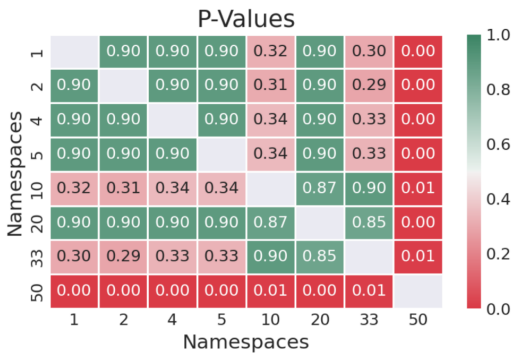
**5.1.2 Ratio of Acceptability.** To augment the Tukey test, a second metric, which we call the Ratio of Acceptability (RoA), is utilized. We define “acceptability” as the range of values for the quantity of interest of the middle 95% of samples in the 1-namespaces case, as seen in Figure 2. If the quantity of interest of an experimental run falls into this range, the replicate is deemed acceptable. The RoA is then defined as the number of samples that are acceptable divided by the total number of samples. The RoA goes beyond the mean comparison of the Tukey test: it indicates the percentage of samples falling within a central 95% region of the benchmark distribution.



**Figure 2:** This plot shows how the upper and lower limit for the range of acceptable values is calculated. The shaded area under the curve represents the middle 95% of alert times for the 1-namespaces case. (a) shows the lower bound of this range, and (b) represents the upper bound.



**(a) Alert time distribution at each namespace.**



**(b) Tukey multiple comparison.**

**Figure 1: Overview of all scanning and detection emulation runs.**

## 5.2 Developing Telemetry Metrics

In addition to establishing a verification process, the second goal of this study is to identify telemetry metrics that can serve as indicators of unreliability for an emulation experiment, which involves selecting telemetry that can be extracted from the emulation and thresholds that the telemetry should not cross for the duration of a given experiment. If thresholds are crossed, then we believe the results of the experiment to be unreliable. Determining a good metric can be done one of two ways: from subject matter expert (SME) input, or through extensive experimentation; in this study we use experimentation to help us determine good indicators. We define a range of acceptable experimental results determined by some baseline set of experiments, and we want to find a telemetry metric and threshold such that the set of experiments that cross this threshold and the set of experiments with results outside the acceptable range are highly correlated. Once we identify telemetry metrics, we can use them to find experiment results that are likely unreliable in future experiments. Although an ideal indicator would be useful across different experiment scenarios, the best telemetry metric and threshold may be highly scenario-dependent.

A good indicator will filter out many unreliable experiments, so that the Tukey and RoA metrics will show better agreement to the baseline compared to the set of results when there was no filtering applied. However, it is also important to look for an indicator that doesn't over-filter the experiments. Thus, the success of different indicators can be compared using the Tukey and RoA metrics described above as well as with a confusion matrix applied to the indicator's results.

The confusion matrix compares the “true” validity of the replicates to the results of filtering with a particular indicator. Validity is determined using the RoA test described above. If the quantity of interest for a given replicate falls within the range of acceptable values given by the RoA test, then the true label for that replicate is “Valid”. Otherwise, the label is “Invalid”. The test label is assigned to the replicate based on the telemetry and threshold for the indicator in question. The confusion matrices are then normalized over the true conditions. A high ratio of false negatives indicates that the telemetry metric is over-filtering reliable replicates. A high ratio of false positives indicates that the telemetry metric is not successfully filtering out replicates that are invalid.

## 6 RESULTS

### 6.1 Scanning and Detection Results

Recall that the method used for oversubscribing resources in this study is to run more and more copies of the experiment scenario in parallel. Experiments for the Scanning and Detection scenario were run under several different emulation configurations, each of which differed in the number of parallel namespaces used to run replicates of the scenario experiment. A summary of these settings can be seen in Table 1. Note that several replicates were removed from the study due to corrupted or insufficient data. Across all settings, there were 3038 total runs with sufficient data for further analysis.

**Table 1: Summary of Scanning and Detection Emulation Settings**

Parallel Namespaces	1	2	4	5	10	20	33	50
Replicates per Namespace	400	200	100	80	40	20	12	8

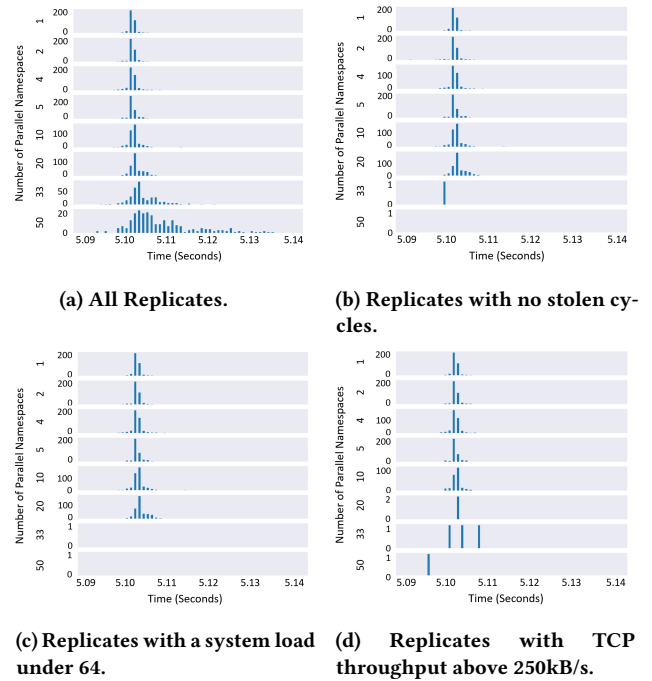
In the scanning and detection scenario, three candidate telemetry metrics were considered independently:

- *Stolen Cycles* – The amount of time stolen from the virtual machine waiting for the host CPU to become available. The threshold for this metric is no cycles stolen during the experiment.
- *System Load* – This is the CPU demand on the physical host in terms of the number of processes running. The threshold for this metric is that the system load will not exceed the number of logical host cores for the duration of the experiment (in this case, 64 cores). Note that system load is calculated as a minute average, and these experiments last on the order of seconds.
- *Throughput* – This is the number of bytes per second of TCP traffic averaged over the course of the experiment. Because the physical host must manage the the TCP traffic in the experiment, this value will change depending on the exact scenario. 250 KB/s, the average of all 1-namespace experimental runs, was chosen as the minimum threshold for this study.

These metrics are captured every second for the virtual machine statistics, and every 10 seconds for the host statistics.

Recall that we consider the 1-namespace experiments the baseline set. In Figure 3(a), we can see the distribution of data deviates significantly from the baseline for numbers of parallel namespaces greater than 20. From these analyses, we would say that the results from 20-, 33-, and 50- parallel namespaces come from different distributions than the baseline data, and that at least some of the replicates under these conditions are unreliable. The other histograms in Figure 3 show the data distributions when certain replicates are filtered out according to the specified indicator. All three metrics appear to do a reasonable job of removing replicates with alert times dissimilar to the baseline distribution.

However, looking at the Tukey comparisons and confusion matrices in Figure 4, the indicators are giving ideal results. System load and throughput do the best job of leaving replicates with alert times that match with the 1-namespace distribution. Stolen cycles



**Figure 3: Histograms of the distribution of alert times at each namespace.**

is actually not as successful as a telemetry metric, according to the Tukey results. However, the Load-based indicator does have a higher false positive rate and the Throughput-based indicator has a higher false negative rate. In fact, in all cases, each of these metrics seem to filter out many replicates that have acceptable alert times – especially those at higher namespaces. This can be seen in Figure 5.

Overall, these results could indicate that the thresholds for the telemetry could be better tuned, although this may make the indicators highly specific to this particular scenario and emulation platform configuration. For all indicators, although the Tukey and confusion matrix results are not perfect, there is a clear improvement in agreement to the baseline distributions when replicates are filtered compared to the unfiltered replicate set. Performance on this specific scenario should be weighed against the desire for generalizable indicators.

The results of the Scanning and Detection study show considerable success in applying the proposed verification procedure to an emulated scenario. We are clearly able to identify the point at which physical resources are over-subscribed. Based on this, we can identify several promising telemetry-based indicators that could be used to help filter out experiments with unreliable results.

Now that we have seen success on one scenario, we aim to apply the same process to the Command and Control scenario and see how the results differ.

### 6.2 Command and Control Results

As with the Scanning and Detection scenario, Command and Control experiments were run under different emulation configurations,

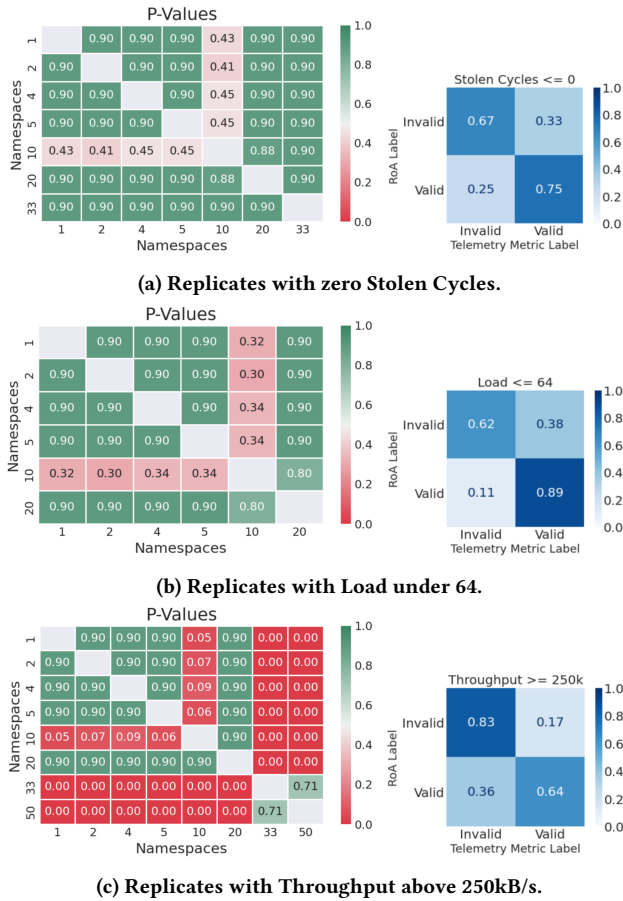


Figure 4: Tukey and Confusion Matrix Results for Selected Indicators on Scanning and Detection Scenario.

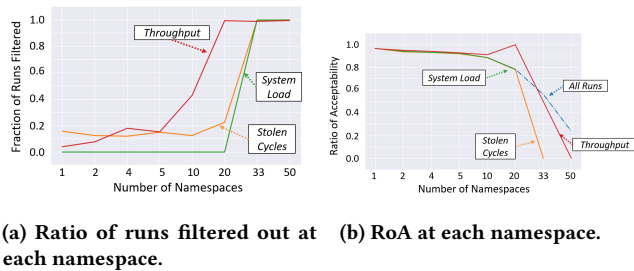


Figure 5: Comparing the number of filtered replicates to the improvement in the RoA metric.

each with increasingly more parallel namespaces. A summary of the settings tested for the Command and Control scenario can be seen in Table 2. The total number of replicates for the command and control verification study is 1187.

Recall that the quantities of interest for this scenario are the number of alerts by time  $t=1, 5, 10,$  and  $16$  seconds. Due to the multiple quantities of interest, Tukey statistics must be aggregated across

Table 2: Summary of Command and Control Emulation Settings

Parallel Namespaces	1	2	5	10	20	40
Replicates per Namespace	200	100	40	20	10	5

all quantities of interest in order to compare result distributions. The aggregate function used in this study was the mean.

Figure 6 shows the Tukey analysis performed on all replicates for this scenario. We can see that around 20 parallel namespaces, the quantity of interest results deviate significantly from the baseline. This is also reflected in the histograms for the scenario. From these analyses, we would say that some of the replicates from 20- and 40-parallel namespaces are likely unreliable.

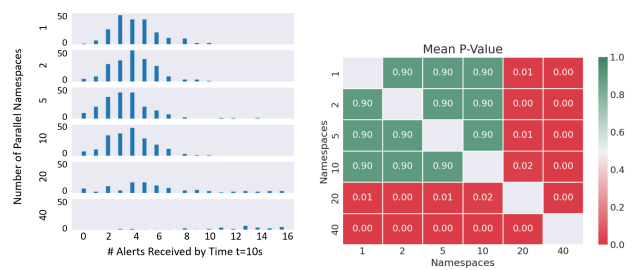


Figure 6: Overview of all command and control emulation runs.

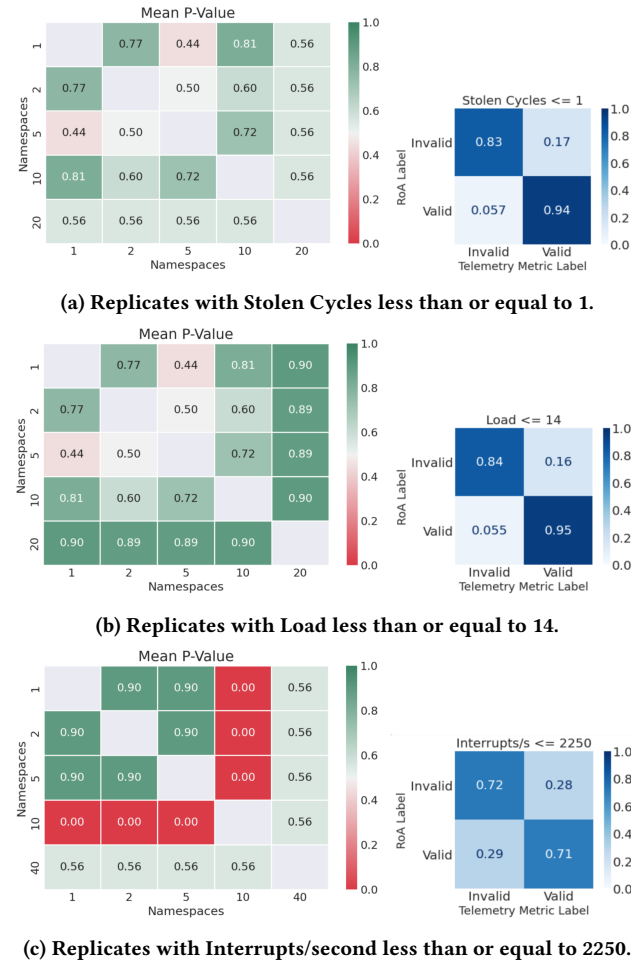
Figure 6: Overview of all command and control emulation runs.

We wished to first test the exact same indicators as the Scanning and Detection scenario. If the same metrics worked approximately as well in both scenarios, it could indicate a highly generalizable metric. In addition, we tested a variety of other available telemetry and thresholds. During experimentation, telemetry is captured every 5 seconds for the VM statistics and every 10 seconds for the host statistics. The set included:

- *Stolen Cycles* – See Scanning and Detection results for telemetry description. The best-performing threshold tested was no more than one stolen cycle.
- *System Load* – See Scanning and Detection results for telemetry description. Although we tested a threshold matching the number of available cores (as with the Scanning and Detection scenario), the Command and Control scenario performed better with a lower threshold, so we chose 14 processes.
- *Interrupts per Second* – This is the number of times in one second that the system is interrupted. The more interrupts there are, the more clock time a single process will take to complete. From looking at the resulting data, a threshold of no more than 2250 was set.

Note that although we wanted to test the Throughput metric with this scenario as well, throughput data was not successfully collected from enough experiments to utilize this metric.

The Tukey comparisons and associated confusion matrices for these indicators can be seen in Figure 7.



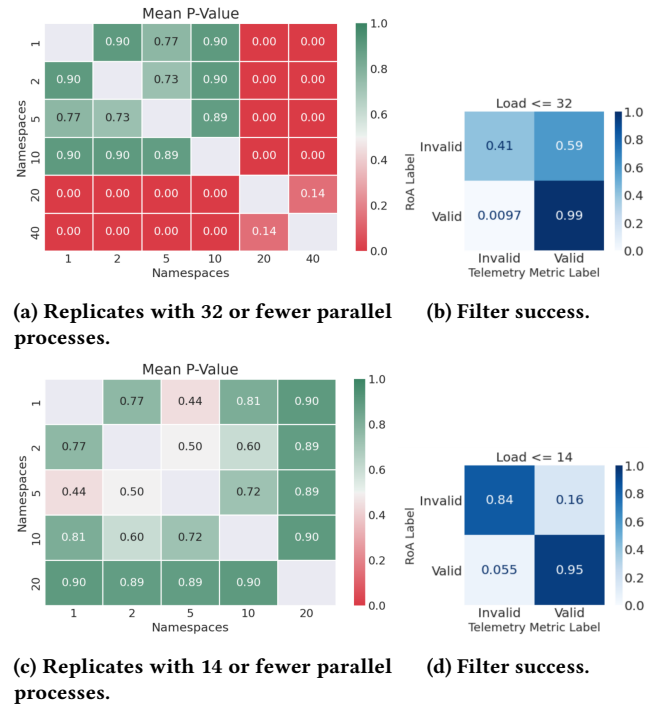
**Figure 7: Tukey and Confusion Matrix Results for Selected Indicators on Command and Control Scenario**

The approach to verification analysis taken with the Scanning and Detection scenario needed to be adapted for the Command and Control scenario. Even when the same telemetry metric was used, the filtering thresholds often needed to be adjusted for the new scenario (6.2.1). Although there isn't an obvious candidate for a universal indicator between the scenarios, we were able to find good indicators for the Command and Control scenario specifically.

**6.2.1 Setting Thresholds.** The threshold for each indicator should be tuned to balance improving the Tukey results overall while also avoiding over-filtering, allowing us to maintain as many reliable replicates as possible.

We first compared the Stolen Cycles indicator with a threshold of 0 (used in the Scanning and Detection Scenario) vs a threshold of 1. When using a threshold of 1, the overall Tukey results did slightly improve, but more importantly, the more forgiving threshold doesn't filter out reliable replicates unnecessarily.

Figure 8 compares the host load indicator with different thresholds. For the Scanning and Detection scenario, the threshold was 64 to match the number of host cores. For reasons discussed in subsection (6.3), the number of cores available on the host was reduced to 32 for the Command and Control scenario. Therefore, the load-based indicator was tested with threshold 32, but an even lower threshold of 14 was found more effective at filtering out the unreliable replicates.



**Figure 8: Tukey and confusion matrices when using different thresholds on "load" telemetry.**

Additional insight can be gained by looking at the distribution of telemetry values in each emulation setting. This is one tactic for approximating thresholds for our metrics, as long as we accept that the threshold will be scenario-dependent. Alternately, such data can also identify telemetry that would make an ineffective indicator.

Figure 9 shows two box plots, one for load and one for number of interrupts per second. We see a smooth increase in load as the number of parallel namespaces increases, with a sharper increase at 20 and 40 parallel namespaces. We can conceive of a horizontal line that might linearly separate datapoints of "good" replicates from "bad" replicates. This separability indicates that load could be a good indicator, and the plot gives us an idea of how to set the threshold. On the other hand, notice the wide range of collected values for the interrupts telemetry. The box plots are mostly centered around the same values, but there are many outliers. This pattern holds regardless of the number of parallel namespaces. These factors could indicate that developing an indicator based on interrupts would be difficult or unsuccessful.

Threshold identification is challenging. Ideally, we want to identify thresholds which are not scenario-dependent, especially because one may not be able to generate controlled conditions in larger, more complicated emulation scenarios. Threshold setting for verification testing deserves more attention in the emulation community.

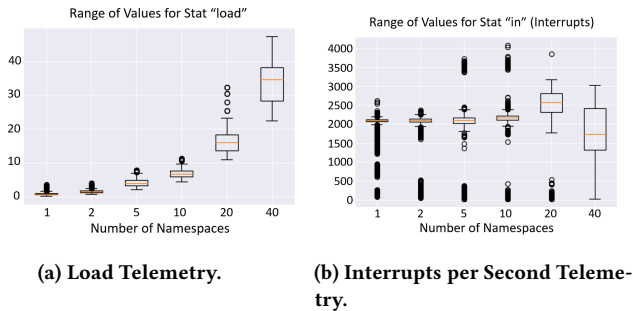


Figure 9: Boxplot of Data Collected Across All Replicates.

### 6.3 Process Generalizability

The verification process we established showed considerable promise on the Scanning and Detection scenario. However, the process was not immediately applicable to the Command and Control scenario. Additional factors had to be considered, such as the difference in resources required to run each scenario, the amount of data collected from each scenario, and how system modeling decisions affect our verification methods.

The Scanning and Detection results (Figure 1) show a clear point where the physical host was oversubscribed, around 33 parallel namespaces. We did not initially see this same point in the Command and Control results - all distributions looked approximately the same regardless of the number of parallel namespaces run. In addition the Tukey analysis showed no p-values lower than 0.05, so we could not reject the hypothesis that all experiments were reliable.

There are a couple factors we theorized would cause the Command and Control results to be inconclusive. The Command and Control scenario was simplified during implementation. The Scanning and Detection scenario involved emulating 27 VMs while the Command and Control scenario only contained four. The Scanning and Detection scenario had one host connecting to many, which required TCP connection maintenance throughout the experiment. The Command and Control scenario simplified network traffic emulation by having only one VM output traffic. Thus, the Command and Control scenario required far fewer resources to run than the Scanning and Detection scenario, so the physical host was not as easily oversubscribed.

One approach to address this was to run more parallel namespaces. However, we found during testing that our emulation infrastructure would not run 50 or more parallel namespaces. Contact authors for further discussion. Another approach to increasing resource demand is to reduce the amount of resources available to the host. We did so by reducing the available processing power on

the physical host by half. The results described in 6.2 came from a set of experiments where resources were limited in this way. From the results shown, it is clear that this configuration did lead to oversubscription of resources, as desired. Several additional approaches are left to future work (Section 7).

With the conclusion of the two scenario studies, we have shown that our proposed verification process is generalizable and gained a great deal of insight.

- More consideration was required to process Command and Control data in a timely and meaningful manner due to the volume of data. Modifications to data collection methods were required to preemptively reduce the volume of output traffic data, and even then 1200 experiment replicates output approximately 20Gb of data total.
- Verification can highlight subtleties in physical host configurations. Using our local computing cluster, we found a marked difference in results depending on which physical host ran the experiments, even though each host should have similar specifications. This is not an unprecedented finding, as previous work by Maricq et al. [20] noted similar performance variability.
- The design of the emulated system model can effect verification methods and results. When addressing the level of fidelity with which to model a system, consider the abilities of the physical host running the emulation. Similarly, when telemetry is collected, consider the emulation model when interpreting results.

Overall, the verification process we have outlined presents a significant step forward in emulation verification. We were able to use our process to verify the capabilities of our physical host to run large numbers of parallel emulations successfully without introducing emulation artifacts to the resulting data. In addition, our process shows great promise of being a generalizable framework for emulation verification.

## 7 FUTURE WORK

We have only scratched the surface of emulation verification. First, there are a couple fundamental assumptions on which our verification analysis relied. We made sensible assumptions at the time to build this study, but it may be fruitful to look back and reconsider the following:

- The acceptable distribution of experiment results was determined based on a baseline set of data. This baseline was calculated from results of experiments run one at a time, minimizing resource utilization. But if the scenario was sufficiently large and complex, or the physical resources were sufficiently limited, only a single experiment may overwhelm the physical host, giving unreliable results. The assumption that the physical host has enough resources for one experiment was needed to help “prove” our hypotheses, but proving this assumption is true could be challenging.
- We define “sameness” of distributions by a successful Tukey test. But is this enough? How close do two distributions have to be to each other to be considered the same? And is this sameness of distributions sufficient to indicate reliable experiment results?

Second, we have only explored one definition of verification: resource verification. There are other levels of realism to verify, including network traffic patterns and software execution[15]. Different experimental goals may require different types of verification, and in turn require different metrics.

Third, the Results section (6) discussed two methods we pursued to force the Command and Control scenario to stress the available physical resources of the system. An alternate approach is to make the scenario itself more complex:

- Add more VMs to further stress memory resources.
- Have the VMs perform more resource-intensive processes.
- Make the network for the scenario more complex, stressing the host by forcing it to maintain more emulated network connections.

Performing these experiments could lead to new research questions, such as which tactics are more effective at causing different types of stress to the system, which could add new context to telemetry metrics.

Finally, there is much still to explore about indicator formulation. For example, consider the role of machine learning in indicator development with subject matter expert opinion. Models could weight and consider multiple telemetry sources, and multi-telemetry metrics could be a more comprehensive evaluation of experiment health and reliability as compared to a single-telemetry metric. In addition, high feature importance in the model could also drive further single-telemetry indicator studies. There may be many other approaches to developing informative indicators of reliability.

Note that our experimental setup may be replicated using the minimega [22] platform. While the SCORCH tool [14] is not currently open source, the main purpose of SCORCH in this study was to facilitate running multiple experiments in parallel, which may be achieved through other means. The scenarios used in this study can be recreated from the descriptions in [30], [29], and [31], but it would also be informative to explore the described verification process and avenues of future work in the context of new scenarios.

## 8 CONCLUSIONS

Verification is key to ensuring that system emulation successfully models the real world. We demonstrated our proposed verification process on two different cyberattack scenarios. In the Scanning and Detection scenario, we found that we could use telemetry-based metrics to detect experiments that were likely to have unreliable results. By applying a similar method to the Command and Control scenario, we were able to show that the same verification process, while not immediately applicable, was fairly generalizable to the new scenario. The exercise also gave us several key insights into important considerations for verification.

In general, experiment repeatability is key to bringing rigor and scientific process to cybersecurity. Just as reported results in the physical sciences should be repeatable and verifiable, so too should cyber experiments. Failure to reproduce a result may be the fault of emulation artifacts rather than faulty experiment design. The process of verification and use of telemetry metrics can help evaluate these factors. Reproducibility of experiments strengthens belief in the validity of results. A result that holds under repetition and various modeling environments is more likely to hold in the real

system as well. In this way, verification can give us more confidence that experimental environments can inform our decisions on real life cyber systems.

## ACKNOWLEDGMENTS

This work has been supported by the SECURE project within the Grand Challenge LDRD Program at the Sandia National Laboratories. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology & Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525. This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

## REFERENCES

- [1] 2012. IEEE Standard for System and Software Verification and Validation. *IEEE Std 1012-2012 (Revision of IEEE Std 1012-2004)* (2012), 1–223. <https://doi.org/10.1109/IEEESTD.2012.6204026>
- [2] 2018. *Alert (TA18-201A) Emotet Malware*. Technical Report. US CERT. <https://us-cert.cisa.gov/ncas/alerts/TA18-201A>
- [3] 2021. Azure Monitor. (2021). <https://docs.microsoft.com/en-us/azure/azure-monitor/autoscale/autoscale-common-metrics>.
- [4] 2021. Cloud Telemetry: Advancing Your IT Strategy. (2021). <https://www.intel.com/content/www/us/en/cloud-computing/telemetry.html>.
- [5] 2021. Network Telemetry. (2021). <https://cloud.google.com/network-telemetry>.
- [6] 2021. What is Telemetry? The Guide to Application Monitoring. (2021). <https://www.sumologic.com/insight/what-is-telemetry/>.
- [7] Terry Benzel, Bob Braden, Ted Faber, Jelena Mirkovic, Steve Schwab, Karen Sollins, and John Wroclawski. 2009. Current developments in DETER cybersecurity testbed technology. In *2009 Cybersecurity Applications & Technology Conference for Homeland Security*. IEEE, 57–70.
- [8] Scott Brown, Brian Henz, Harold Brown, Michael Edwards, Michael Russell, and Jonathan Mercurio. 2015. Validation of network simulation model with emulation using example malware. In *2015 IEEE Military Communications Conference (MILCOM)*. IEEE.
- [9] Prasad Calyam, Sudharsan Rajagopalan, Sripriya Seetharam, Arunprasad Selvadurai, Khaled Salah, and Rajiv Rammath. 2014. VDC-Analyst: Design and verification of virtual desktop cloud resource allocations. *Computer Networks* 68 (2014), 110–122. <https://doi.org/10.1016/j.comnet.2014.02.022> Communications and Networking in the Cloud.
- [10] Cisco. 2019. Snort intrusion detection and prevention system. <https://www.snort.org/>.
- [11] Jonathan Crussell, Thomas M Kroeger, Aaron Brown, and Cynthia Phillips. 2019. Virtually the Same: Comparing Physical and Virtual Testbeds. In *2019 International Conference on Computing, Networking and Communications (ICNC)*. IEEE.
- [12] Jon Davis and Shane Magrath. 2013. *A survey of cyber ranges and testbeds*. Technical Report. Cyber and Electronic Warfare Division, Defence Science and Technology Organization, Australian Government.
- [13] Bernard Ferguson, Anne Tall, and Denise Olsen. 2014. National cyber range overview. In *2014 IEEE Military Communications Conference*. IEEE, 123–128.
- [14] Seth Hanson, Jerry Cruzy, and Casey Glatter. 2021. *SCORCH User Guide*. Technical Report SAND2021-11504 O. Sandia National Laboratories.
- [15] Brandon David Heller. 2013. Reproducible Network Research with High-Fidelity Emulation. (2013).
- [16] Mike Hibler, Robert Ricci, Leigh Stoller, Jonathon Duerig, Shashi Guruprasad, Tim Stack, Kirk Webb, and Jay Lepreau. 2008. Large-scale Virtualization in the Emulab Network Testbed. In *2008 USENIX Annual Technical Conference (USENIX ATC 08)*. USENIX Association.
- [17] Alefiya Hussain, David DeAngelis, Erik Kline, and Stephen Schwab. 2020. Replicated Testbed Experiments for the Evaluation of a Wide-range of DDoS Defenses. In *2020 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*. IEEE, 46–55.
- [18] Quan Jia, Zhaohui Wang, and Angelos Stavrou. 2009. The Heisenberg Measuring Uncertainty in Lightweight Virtualization Testbeds. In *CSET*.
- [19] Junyu Lai, Jiaqi Tian, Ke Zhang, Zheng Yang, and Dingde Jiang. 2021. Network Emulation as a Service (NEaaS): Towards a Cloud-Based Network Emulation Platform. *Mobile Networks and Applications* 26 (2021), 766–780. Issue 2. <https://doi.org/10.1007/s11036-019-01426-0>

- [20] Aleksander Maricq, Dmitry Duplyakin, Ivo Jimenez, Carlos Maltzahn, Ryan Stutsman, and Robert Ricci. 2018. Taming Performance Variability. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI'18)*. USENIX Association.
- [21] Stephen Midway, Matthew Robertson, Shane Flinn, and Michael Kaller. 2020. Comparing multiple comparisons: practical guidance for choosing the best multiple comparisons test. *PeerJ* 8 (2020). <https://doi.org/10.7717/peerj.10387>
- [22] minimega developers. 2019. minimega: a distributed VM management tool. <http://minimega.org/>
- [23] Jelena Mirkovic, Genevieve Bartlett, and Jim Blythe. 2018. DEW: Distributed Experiment Workflows. In *11th USENIX Workshop on Cyber Security Experimentation and Test CSET 18*.
- [24] Jelena Mirkovic, Terry V Benzel, Ted Faber, Robert Braden, John T Wroclawski, and Stephen Schwab. 2010. The DETER project: Advancing the science of cyber security experimentation and test. In *2010 IEEE International Conference on Technologies for Homeland Security (HST)*. IEEE, 1–7.
- [25] William L Oberkamp and Christopher J Roy. 2010. *Verification and validation in scientific computing*. Cambridge University Press.
- [26] Lee M Rossey, Robert K Cunningham, David J Fried, Jesse C Rabek, Richard P Lippmann, Joshua W Haines, and Marc A Zissman. 2002. Lariat: Lincoln adaptable real-time information assurance testbed. In *Proceedings, IEEE Aerospace Conference*, Vol. 6. IEEE, 6–6.
- [27] Stephen Schwab and Erik Kline. 2019. Cybersecurity Experimentation at Program Scale: Guidelines and Principles for Future Testbeds. In *2019 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*. IEEE, 94–102.
- [28] Christos Siaterlis, Andres Perez Garcia, and Béla Genge. 2012. On the use of Emulab testbeds for scientifically rigorous experiments. *IEEE Communications Surveys & Tutorials* 15, 2 (2012), 929–942.
- [29] Thomas Tarman, Trevor Rollins, Laura Swiler, Jerry Cruz, Eric Vugrin, Hao Huang, Abhijeet Sahu, Patrick Wlazlo, Ana Goulart, and Kate Davis. 2021. Comparing reproduced cyber experimentation studies across different emulation testbeds. *Proceedings, 14th Workshop on Cyber Security Experimentation and Test (CSET21)*. ACM (2021).
- [30] Eric Vugrin, Jerry Cruz, Christian Reedy, Thomas Tarman, and Ali Pinar. 2020. Cyber threat modeling and validation: port scanning and detection. In *Proceedings of the 7th Symposium on Hot Topics in the Science of Security*. Association for Computing Machinery.
- [31] Eric Vugrin, Seth Hanson, Jerry Cruz, Casey Glatter, Thomas Tarman, and Ali Pinar. 2021. Detection of command and control traffic: model development and experimental validation. *in preparation* (2021).
- [32] Xi Zheng and Christine Julien. 2015. Verification and Validation in Cyber Physical Systems: Research Challenges and a Way Forward. In *2015 IEEE/ACM 1st International Workshop on Software Engineering for Smart Cyber-Physical Systems*. IEEE. <https://doi.org/10.1109/SEsCPS.2015.11>