

Towards A Formal Framework for Evaluating the Effectiveness of System Diversity when Applied to Security

Raphaël Khoury*, Abdelwahab Hamou-Lhadj[†] and, Mario Couture[‡]

*Defence Research and Development Canada

Valcartier, Quebec, Canada,
raphael.khoury@drdc-rddc.gc.ca

[†]Concordia University

Software Behaviour Analysis (SBA) Research Lab
Department of Electrical & Computer Engineering
Montreal, Quebec, Canada
abdelw@ece.concordia.ca

[‡]Defence Research and Development Canada

Valcartier, Quebec, Canada,
mario.couture@drdc-rddc.gc.ca

Abstract—N-version programming has been shown to be an effective way to increase the reliability of systems. In this study, we examine the possibility of extending this approach to address security, rather than reliability concerns. We focus specifically on how to evaluate the efficiency of the use of diversity for security. We show that while several key elements must be taken into account when N-version programming is used for security rather than reliability, it is nonetheless possible to devise a reasoning framework to evaluate the efficiency of this development paradigm in a security context. This framework allows us to reason about the most effective way to use diversity for security.

I. INTRODUCTION

The fields of software reliability and security are closely related and several methods simultaneously address both concerns without distinguishing between a malicious and an inadvertent failure. It is thus normal to ask if the N-version programming paradigm, which was developed to address reliability concerns, can likewise be deployed in a security context. We believe this is possible, but that several key elements must be taken into account when diversity is introduced in an architecture for security purposes rather than to increase reliability.

N-version programming [1] is a software development paradigm that draws upon the concept of diversity to increase the reliability of software. The guiding principle of this approach is to produce several distinct versions of a given software, and execute them in parallel with the same inputs. A discrepancy between the outputs of the various instances is an indication that at least one instance has malfunctioned. In that case, a single output value is chosen from the outputs of each program instance by majority voting. The intuition behind this programming paradigm is that while it may be impossible to produce a single flawless instance of any complex system,

multiple instances of this system would normally exhibit different faults.

A key concept in the design of N-version architectures is *failure independence*. Informally, this property describes the behavior of a system for which the occurrence of a failure in one instance for a given input value does not provide any information in regard to the probability of failure of another instance for the same input value. It is from the assumption of failure independence that we derive the hypothesis that the probability of *coincident failure* (i.e. two instances failing on the same input) is very small and that gains in reliability can be obtained through the use of N-version programming.

Researchers in security have also shown a great deal of interest in diversity, though not in the context of an N-version architecture. Instead, research in computer security proceeds from the assumption that if the program instance of each user differed from that of every other user, an attack cannot easily be carried over from one system to the next. The attacker will thus be forced to tailor each attack to the system he wishes to compromise. Furthermore, the added uncertainty about the target system increases the cost of the attack [2].

Diversity can thus serve as the basis for effective intrusion detection. The main intuition underlying such an approach is that, since attacks must be tailored to each program instance, if several program instances are executed in parallel as part of an N-version architecture and an input contains an attack vector, it is likely that the attack will succeed only on some of the several program instances. This in turn will cause the executions of the affected and unaffected instances to diverge observably from one another. Such a divergence can then serve as the basis for intrusion detection and reaction. Initial research shows this approach to be highly promising [3].

In this study, we propose a framework that can be used to

evaluate the effectiveness of the use of diversity for security. We contrast the use of diversity for security purposes to that of diversity for reliability and highlight a number of key differences that must be taken into consideration in the former case. We show that while several key elements must be taken into account when N-version programming is used for security rather than reliability, it is nonetheless possible to devise a reasoning framework to evaluate the efficiency of this development paradigm in a security context.

The use of our proposed framework would bring several benefits. Most notably:

- This framework allows us to reason about the impact of diversity on the security of a system, and thus determine the most effective way to introduce diversity for security purposes.
- It allows us to compare architectures in which diversity is introduced at different layers of the system (OS, hardware or software), or to compare the effectiveness of diversified architectures built from different combinations of COTS software.
- It can provide indications as to which architecture we can expect to provide the best detection rate or the lowest rate of false positives, and thus serve as a guide for design choices.

The remainder of this paper is organized as follows. In Section II we contrast the proposed use of N-version programming for security with its more common use for reliability and highlight the relevant differences between the two approaches. Section III reviews the literature on both topics. Concluding remarks and perspectives for future work are presented in Section IV.

II. DIVERSITY FOR SECURITY VS DIVERSITY FOR RELIABILITY

A. General Framework

We propose the following framework to study and reflect about the use of diversity for security. We start with a population \mathcal{P} of programs, that represents a hypothetical set of all possible programs able to solve a given problem. We let π range over possible programs.

$$\mathcal{P} = \{\pi_1, \pi_2, \pi_3, \dots\}$$

These programs take their input from a set of possible input values \mathcal{X} . Each input represents the entire interaction a user has with a given program during a session. We let x range over the possible input values.

$$\mathcal{X} = \{x_1, x_2, x_3, \dots\}$$

Some inputs may be *malicious*, meaning that they hide exploits that bring the system in a state that violates the security policy. For instance, an input field may contain data triggering a buffer overflow and allowing code injection to occur. Such input values are said to be *invalid*. Normal input values that do not contain an attack, are said to be *valid*. We write \mathcal{X}_v for the set of valid inputs and \mathcal{X}_i for that of invalid inputs. Note that $\mathcal{X} = \mathcal{X}_v \cup \mathcal{X}_i$ and that \mathcal{X}_v and \mathcal{X}_i are disjoint.

Let $Q(\cdot)$ stand for the usage distribution of the values of \mathcal{X} . This distribution naturally affects only valid executions, since

an attacker can alter the distribution of inputs by repeatedly inputting the values he or she needs in order to alter the system.

In the context of their study of diversity for security, Littlewood et al. [4] propose a score function $v : (\mathcal{P} \times \mathcal{X}) \rightarrow \{1, 0\}$, that indicates whether or not the execution of a given program for a given value fails. The occurrence of such a failure could be immediately observed by inspecting the program's output (for example if the program failed to produce a return value), or discovered by contrasting this output with that of another instance. The score function v thus serves as basis for evaluating how reliable a given program is and a given diverse architecture containing this program.

However, when the focus is on security, the program's output alone does not provide sufficient information for a meaningful evaluation of the validity of the input. It is entirely possible for an intruder to alter the execution in such a way as to violate the security policy while keeping the output unchanged. Indeed, an attacker who wishes to remain undetected would favor such a course of action and many attacks are not visible at the level of the output alone. It follows that a diversity-based framework whose focus is security rather than reliability, will necessarily rely on a complete trace of the program's execution, rather than simply on the output, as the basis for its evaluation of the input. This is the first main difference between using diversity for reliability and using it for security.

Difference 1: *Diversity for reliability is implemented by comparing the outputs of multiple instances. When diversity is used for security purposes however, it is necessary to examine execution traces.*

An execution trace is a sequence of atomic actions, performed by the target program during its execution and recorded by a reference monitor. It can hypothetically contain every instruction performed by the target program, or be focused on a subset of security-relevant actions tailored to the security policy of interest or to a specific resource whose security we seek to optimize. Let Σ stand for the set of all possible execution traces and let σ range over traces.

The trace function $\nu : (\mathcal{P} \times \mathcal{X}) \rightarrow \Sigma$ thus replaces the score function of Littlewood et al. This function is given as :

$$\nu(\pi, x) = \sigma \text{ where } \sigma \text{ represents the trace of program } \pi \text{ on input } x. \quad (1)$$

Once the system is executed on multiple instances, it is necessary to contrast the execution traces in order to detect a possible intrusion. Let $\sigma_1 = \nu(\pi_1, x)$ and $\sigma_2 = \nu(\pi_2, x)$ be the execution traces of programs π_1 and π_2 respectively over the same input value x . Let the correlation function $corr(\sigma_1, \sigma_2)$ stand for the degree of similarity observed between these two executions. This correlation is expressed by way of a value between 0 and 1, where 1 indicates identical sequences, and 0 identifies sequences that seem completely unrelated. Several

metrics could be used to compute this value. A natural choice is the Levenshtein distance [5], a measure of the number of insertions, deletions and replacements needed to turn one sequence into the other. A similar metric was successfully used in [6] for intrusion detection. However, other metrics specifically designed for the problem of detecting divergence between program executions could also be considered.

$$\text{corr}(\sigma_1, \sigma_2) = \text{The degree of similarity between executions } \sigma_1 \text{ and } \sigma_2. \quad (2)$$

The central idea that underlies the use of diversity as a defence mechanism is that a given attack may succeed on one system but fail on another. This in turn will lead to an observable divergence in the execution traces, allowing the attack to be detected. As discussed above, diversity can be introduced at various levels, such as memory layout [7], instruction set [8] or by using two distinct implementations of the same software or operating system. The success of the approach rests on the capacity to develop systems that are sufficiently different so that most attacks cannot succeed on multiple instances. It follows that while developers building a N-version architecture with the goal of increasing reliability must focus on minimizing the occurrence of coincident failure, those building such architectures for security purposes should seek to maximize the dissimilarity of internal behavior.

Difference 2: *In the context of diversity for reliability, the design of an N-version architecture must minimize the occurrence of coincident failure. However, if the object is security, the design should promote dissimilarity of behavior.*

This second difference raises several interesting questions related to the way to maximize the divergence between systems while maintaining common functionalities between instances, as well as the way to simultaneously update both instances so as to preserve their behavioral equivalence.

Pairs of systems naturally differ as to how much similarity they exhibit while executing normally (i.e. over valid inputs). However, a baseline can be established by examining a sufficiently large and representative sample of executions. This yields a distribution θ as follows:

$$\theta(\pi_1, \pi_2) = \sum_i \text{corr}(\nu(\pi_1, x_i), \nu(\pi_2, x_i))Q(x_i) \quad (3)$$

In effect, this distribution expresses the likelihood that when executing a given input, two executions will differ by a given amount¹. Note that this distribution is only computed for valid executions. We expect that an invalid execution for which the attack succeeds on one instance only, will contrast with a corresponding valid execution by a higher than average amount, but at the present time this can only be a conjecture.

¹Observe that this distribution only describes the behavior of a specific version of each system.

Let $\text{corr}(\sigma_1, \sigma_2)$ be the level of similarity existing between two trace executions σ_1 and σ_2 . Our goal is to contrast the observed $\text{corr}(\sigma_1, \sigma_2)$ with the known value of $\theta(\pi_1, \pi_2)$ of the program that produced σ_1 and σ_2 , as to attempt to determine if the input value x hides an attack. Were we in possession of statistical data about the relative distribution of valid and invalid inputs, as well as of data relating to the expected level of correlation between executions of the target programs over invalid inputs, a statistical analysis could be performed. Such an analysis could return a probability ϕ indicating that x is malicious with a certain degree of confidence. Indeed, this is similar to the method used to evaluate the reliability of diversified components. However, as discussed above, it is not meaningful to compute a distribution of invalid inputs when reasoning about the possible behavior of a malicious attacker capable of altering the systems's inputs. Furthermore, while statistical data could be gathered about systems behavior on malicious inputs by simulating their execution using test cases of known attacks, such data may not necessarily be generalized to new or unknown attacks. Of particular interest are zero-day exploits that use unknown vulnerabilities of software.

We propose instead a possibilistic approach. Possibility theory [9], is an alternative to probability theory to reason about uncertainty. Informally, a possibility is a value between 0 and 1 that describes the ease with which an event will occur, or will belong to a set, as opposed to the likelihood that it will occur, which is expressed as a probability.

A possibilistic analysis would thus indicate how “normal” the level of observed correlation is and how unusual it would be for a valid input to result in pair of executions exhibiting this level of correlation. This information is captured by a possibility function pos whose domain is the range of possible correlation values and whose image is a possibility value in the interval $[0, 1]$. Dubois et al. [10] show how a possibility function reflecting this value can be computed directly from the probability density function. The technique they propose can thus be used to derive a possibility distribution for correlation values. Formally, the possibility of a correlation u occurring is

$$pos(u) = \int_0^u \theta(y)dy + \int_{f(u)}^1 \theta(y)dy \quad (4)$$

where $f : [a, u_0] \rightarrow [u_0, b]$ is a function defined s.t. $f(u) = \max\{y | \theta(y) \leq \theta(u)\}$ with the interval $[a, b]$ being the support of θ (possibility $[0, 1]$) and u_0 being its modal value².

Once a possibilistic value has been assigned to the correlation between two sequences, action can be taken based on the environment-specific tradeoff between security and functionality.

One way to control this tradeoff is to state a threshold α , a level of divergence below which any pair of executions is deemed suspect.

²This solution is thus only applicable to the cases where θ is continuous and monomodal, but it is reasonable to think that this will usually be the case of any correlation probability density function.

$$\vartheta = \begin{cases} 1, & \text{if } \phi \geq \alpha; \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

The choice of the threshold value α will determine the sensitivity of the detection. However, in the absence of data about the level of divergence observable in invalid execution it is not possible to give a numerical value to our confidence in this judgement.

Once the attack is detected, the system administrator can react, usually by terminating the execution. In this distinction lies another consequential difference between using diversity for reliability and using diversity for security: in the first case, the goal is to maximize the number of input values for which a correct service is provided, whereas in the latter, the goal is to weed out and deny service if the input value betrays a malicious intent on the part of its originator. While information from the unaffected instance may be used to recover from the attack, it is undesirable to provide service to a malicious user since doing so betrays information about the system.

Difference 3: *When diversity is implemented with the goal of increasing reliability, the object is to maximize the number of inputs for which a service is provided. In the case of diversity for security, we seek to weed out invalid inputs.*

This is more than a simple difference in the reaction to the discovery of a vulnerability and translates into profound changes as to how an architecture must be evaluated. In particular, it poses a new risk which does not exist when diversity is implemented for reliability: that of the occurrence of false positives when two valid executions diverge to the point that the monitor wrongly marks one of them as being invalid.

Difference 4: *The fact that some inputs will not be answered, coupled with the imperfection present in any security architecture, leads to the occurrence of false positives.*

The likelihood that a pair of executions will be mistakenly marked as malicious is a factor of the decision threshold α used to rule out executions and of the distribution θ . We write $fp(\alpha, \theta)$ for the probability of false positives occurring if the threshold is a set of α and the similarity distribution between valid executions is θ .

This risk is expressed by the following equation:

$$fp(\alpha, \theta) = \int_0^y \theta(u) du \quad (6)$$

where y is the maximal value for which $pos(y) \leq \alpha$.

The occurrence of false positives can at first sight be thought of as analogous to the occurrence of coincident failures in the diversity for reliability context and could thus presumably be studied using the same analytical tools that have already been developed for the latter case. This analogy does however obscure two critical differences. First, coincident failures

are the result of *unwanted* commonalities between different instances, and various strategies are employed to minimize and eliminate these commonalities (see for e.g. [11], [12]). False-positives, however, result from the *desired* divergence between instances and we believe that, in general, the approach grows stronger as these differences increase. Future research should determine if there exists an optimal amount of divergence that provides the best ratio of attack detection to false-positives for a given threshold. In the meantime, abstraction and correlation algorithms should be developed to discern the divergence between executions associated with successful intrusion from those which occur naturally because of differences in the underlying programs.

Secondly, the modeling of false positives also differs from that of coincident failures in that the latter relies upon a distribution over all inputs to assess the rate of coincident failures and contrasts it against that of faults which are successfully tolerated by the N-version architecture. However, while a distribution of *valid* inputs can be constructed and the rate of false positives estimated from it, we can never hope to compute a distribution that includes invalid inputs, since the occurrence of invalid inputs implies the presence of an attacker capable of querying the system with inputs of his choice calibrated for the purpose of his attack. The best we can do in this case is thus to estimate the rate of occurrence of false positives of the system when it is not under attack.

It does seem intuitive that as instances grow more different, so does the benefit of using diversity for security. Indeed, it is more likely that an attack will succeed on only one of two instances if they are very different from one another than if they are alike. However, as executions grow more divergent it will also become more difficult to identify similarities and correlations between valid executions. This in turn could lead to an increase in the rate of false positives. It follows that unless the increase in divergence between instances is matched by a corresponding increase in the sophistication of the correlation algorithm, benefits by increasing the divergence between instances will be offset by an increase in the number of false positives

B. Multiple Instances

Research on diversity for reliability shows that the overall reliability of the system can be improved (up to a point) by increasing the number of instances present in a N -version architecture (i.e. by increasing the value of N) [13]. This is a strategy designed to cope with the unavoidable presence of coincident failures: even if two or more instances fail on the same input, other instances running concurrently may succeed. Such multi-versions architectures take two forms: one-out-of- N systems, which succeeds if at least one instance succeeds, and majority voting systems (or m -out-of- N systems, with $m = \lceil N/2 \rceil$) that succeeds if a majority of the composing instances succeed.

The previous idea can also be carried over to the context of diversity for security. Running multiple instances will increase the odds that a malice will be detected by making it

harder on the attacker to devise an input that can compromise every instance simultaneously. A final key difference between diversity for reliability and diversity for security is that, in the latter case, rather than a rigid choice between one-out-of- N or majority voting a diverse architecture can be devised for any $m < N$, indicating that this attack is deemed to have occurred if at least m instances diverge from the others. As the number of instances deviating from the others increases, the input can be treated as malicious with an increasing level of confidence.

Difference 5: *One-out-of- N and majority voting are no longer the only possible voting paradigms. Instead, a threshold dictating how many instances can deviate from the others before the input is treated as malicious must be chosen in such a way as to balance security concerns with the need to minimize false-positives.*

However, comparing several instances raises a number of difficulties. In particular, the system becomes subject to a variation of the consistent comparison problem raised in [14]. Consider a system with three instances, π_1, π_2 and π_3 which, for a given input value x , produce three traces σ_1, σ_2 and σ_3 respectively. Three pairs of comparisons are possible between these three sequences, namely σ_1 with σ_2 , σ_1 with σ_3 and σ_2 with σ_3 . Let ρ_1, ρ_2 and ρ_3 be the outputs of these three comparisons and let ϑ be the threshold by which we consider that an instance has diverged from the other (and thus that the input is malicious). It may become possible that $|\rho_1 - \rho_2| < \vartheta$, $|\rho_2 - \rho_3| < \vartheta$ but that $|\rho_1 - \rho_3| > \vartheta$.

C. System Health and Self Monitoring Execution

In addition to contrasting the various executions between themselves to detect a violation, each execution can be contrasted with a model of the system's desired behavior to detect if the ongoing execution violates the system's security property. This would have several benefits: first, it would reduce the number of false positives by giving an indication as to whether or not a deviation observed between two execution rallies does correspond to a violation of the security policy. Secondly, it would indicate which of two diverging instances is the one for which the attack has succeeded, thus allowing us to isolate the compromised system and use the healthy one for recovery.

The execution monitoring, like the correlation analysis, would return a probability that the execution under consideration is malicious. This is given by the function $eval : \Sigma \rightarrow [0, 1]$.

$$eval(\sigma) = p \text{ the probability that } \sigma \text{ is malicious.} \quad (7)$$

The evaluation given by $eval$ is then contrasted with that of the correlation analysis. We can state with greater confidence that two diverging executions hide an attack if they not only diverge, but if they also have a high probability of maliciousness according to $eval$. Once again, the decision as to whether or not a specific execution is to be treated as malicious

will be based on a tradeoff between security concerns and the desire to minimize the rate of false positives.

In the presence of a self-diagnostic for each execution sequence, the equations for detecting an attack and to compute the rate of false positives can be stated in a variety of ways, depending on how much weight is given to each type of judgement. An elegant solution is to merge the two judgements of the $eval$ functions into a single value that gives the probability that at least one of the two instances is malicious (according to the self-diagnostic) and then adjust the threshold according to this value. The detection should be more sensitive to divergence between the two instances if the self-diagnostic raises alarms, and conversely more tolerant if no abnormal behavior is detected in either instance.

III. RELATED WORKS

The N-programming development paradigm [1], also called design diversity, grew from the longstanding practice of using multiple redundant components to increase the reliability of safety-critical hardware. This practice has a rich literature dating back to the late 70s, that includes various experiments conducted in academic settings to evaluate the feasibility and efficiency of the approach as well as theoretical enquiries aimed at modeling and reasoning about the behavior of N-version systems.

The latter studies begin with Eckhardt et al. [15], who proposed an initial model to study the impact of coincident failures on the effectiveness of using diversity for security. The authors modeled the occurrence of such failures by an intensity function that represents the propensity of programmers to introduce faults in such a way that failure does not occur independently on some inputs. Building upon their work, Littlewood et al. [4] proposed an alternative model, that includes an important dimension of methodology, to model the impact of different development strategies on the distribution of coincident failures in software. Both models are contrasted and discussed in [16]. A final approach is proposed by Partridge et al. [17] to model the distinction between the different ways several instances may fail, even if they fail over the same inputs.

The question of using N-version programming for security, rather than for reliability, was raised in a number of studies. Littlewood et al. examined the question in [18]. In [19], Bessani, et al. argued in favor of using diversity for security on the basis of the recorded distribution of vulnerabilities in several operating systems. In [20], the various layers where diversity can be inserted are examined from the perspective of maximizing security. The use of design diversity to protect against computer viruses were examined in [21]. The use of diversity for security, termed automated diversity, was first suggested by Forrest et al. in [2], where it was observed that the homogeneity of computer systems constitutes an important vulnerability. Drawing on an analogy to biological systems, Forrest et al. argued that the robustness of systems could be improved if the program instance used by each user differed slightly from that of every other user. As a case study, they

proposed a method for stack layout randomization and showed that it is effective at disrupting a buffer overflow attack. In [3], Gao et al. propose an intrusion detection scheme based on the behavioral distance between two processes.

Following on this line of enquiry, several studies have proposed introducing diversity at different layers of software systems. These include the instruction set [8], [22], [23], address space [7], [24], data space [25], and system calls [26]. A survey of these techniques is provided in [27].

IV. CONCLUSIONS AND PERSPECTIVES FOR FUTURE WORK

This study proposes a new reasoning framework to evaluate the effectiveness of the use of diversity for security purposes. In this regard, we combine the complementary fields of design diversity, which is traditionally used for the purposes of increasing the reliability of systems, and automated diversity, which is focused on security.

Several questions must still be answered before diversity can be effectively used for security purposes.

- As discussed above, the central assumption that underlies the use of diversity for security is that if two systems are sufficiently different, attacks targeted at one system cannot be carried over to the other by using the same input values. One of the most pressing topics for further investigation is to test the validity of this hypothesis. In this context, it is important to recall that researchers in the field of diversity for security had assumed that independent development was a sufficient condition to achieve failure independence, until research by [28] showed this was not necessarily the case.
- A second assumption often made in this context is that the more different two systems are, the more likely it is that an attack will succeed on one instance but fail on the other. If this is the case, it becomes necessary to make a difficult compromise between risking a higher rate of false positives by increasing the amount of diversity between the program instances, or risking a lower detection rate by using more similar systems. Answering this question may lead to the development of new metrics to measure the level of divergence between two systems.
- A related problem lies in determining if there are areas of the system for which a greater benefit is derived by introducing diversity. It may happen that introducing diversity in a few key aspects of a system is sufficient to produce two instances for which few attacks can simultaneously succeed. Since the cost of building diverse instances may be prohibitive, targeting the introduction of diversity to a few key subsystems can allow a more cost-effective enforcement.
- Another possible avenue of future research is to use multiple correlations between instances. Each correlation would be based on its own trace abstraction, and focussed on preventing a specific class of attacks or on protecting a specific key resource. This opens up the possibility of feedback oriented analysis, where the detection of a

possible malice in one of the sequences triggers more scrutiny for all other traces before a final decision can be made about the validity of the trace under observation.

- A final ongoing challenge for future research lies in developing trace analysis and trace abstraction tools that are sufficiently refined to distinguish an attack occurring in a given execution but not in another, from the normal divergence present in a pair of homologous execution traces that are simultaneously executing on diverse systems.

We have already performed initial experiments to implement the ideas proposed in this paper. Initial testing was performed with two diversified html servers, running on the same hardware and operating systems and executing the same requests in parallel. The output of both servers is then contrasted to detect a higher than usual divergence between their behaviors. At this stage we have so far experimented only with normal traces.

These experiment consisted first in a short training phase in which the enforcement mechanism recorded the observed behaviors of both systems when running the same input. This was followed by a test phase in which the enforcement mechanism attempted to predict the expected behavior of one of the system based after having monitored the execution of the other. A deviation between this predicted behavior, and the actual observed behavior would be seen as an indication of an attack.

The results we have obtained are highly encouraging. Indeed, in most cases, the observed difference between both executions was minimal. These results indicate that whenever the execution is valid, the proposed intrusion detection method can successfully detect an attack (i.e., the risk of false positives is not prohibitively high).

The next step is to test the approach with attack traces, and determine the efficiency of the approach using the reasoning framework proposed in the previous sections. For lack of time, we were unable to complete this phase of this experiment. Nonetheless, a survey of the relevant literature indicates that in most cases, an intrusion would incur a greater amount of disturbance in the sequence in system calls than the level (often nil) of the difference we observed between the sequences of each matched pair. It follows that the approach proposed in this section will likely be able to detect intrusions effectively with a minimal false positive rate.

This type of experiment can then be replicated with diversity introduced in other components of the system's architecture (such as the OS or hardware), to gather knowledge about the best way to introduce diversity for security. Other types of analyses, such as principal component analysis, or experiments based on machine learning or data mining techniques could also offer useful insight about the effective use of diversity.

REFERENCES

- [1] A. Avizienis, "The n-version approach to fault-tolerant software," *IEEE Transactions on Software Engineering*, vol. 11, pp. 1491–1501, December 1985.

- [2] S. Forrest, A. Somayaji, and D. H. Ackley, "Building diverse computer systems," in *Proceedings of the Sixth Workshop on Hot Topics in Operating Systems*. IEEE Computer Society Press, 1997, pp. 67–72.
- [3] D. Gao, M. K. Reiter, and D. X. Song, "Behavioral distance measurement using hidden markov models," in *Proceedings of the Recent Advances in Intrusion Detection, 9th International Symposium, (RAID) 2006, Hamburg, Germany*, ser. Lecture Notes in Computer Science, D. Zamboni and C. Krügel, Eds. Springer, September 2006, pp. 19–40.
- [4] B. Littlewood and D. R. Miller, "Conceptual modeling of coincident failures in multiversion software," *IEEE Transactions on Software Engineering*, vol. 15, no. 12, pp. 1596–1614, 1989.
- [5] V. I. Levenshtein, "Binary codes capable of correcting deletions, insertions and reversals," *Doklady Akademii Nauk SSSR*, vol. 163, no. 4, pp. 845–848, 1965, original in Russian – translation in Soviet Physics Doklady 10(8):707-710, 1966.
- [6] S. A. Hofmeyr, S. Forrest, and A. Somayaji, "Intrusion detection using sequences of system calls," *Journal of Computer Security*, vol. 6, no. 3, pp. 151–180, 1998.
- [7] H. Shacham, M. Page, B. Pfaff, E.-J. Goh, N. Modadugu, and D. Boneh, "On the effectiveness of address-space randomization," in *Proceedings of the ACM Computer and Communications Security (CCS) 2004*, B. Pfitzmann and P. Liu, Eds. ACM Press, 2004, pp. 298–307.
- [8] G. S. Kc, "Countering code-injection attacks with instruction-set randomization," in *Proceedings of the ACM Computer and Communications Security (CCS) Conference*. ACM Press, 2003, pp. 272–280.
- [9] L. A. Zadeh, "Fuzzy sets as a basis for a theory of possibility," *Fuzzy Sets and Systems*, vol. 100, pp. 9–34, April 1999.
- [10] D. Dubois, H. Prade, and S. Sandri, "On possibility/probability transformations," in *Proceedings of Fourth IFSA Conference*. Kluwer Academic Publishing, 1993, pp. 103–112.
- [11] A. Avižienis, M. R. Lyu, and W. Schutz, "In search of effective diversity: A six-language study of fault-tolerant flight control software," in *Proceedings the IEEE Eighteenth Annual International Symposium on Fault-Tolerant Computing (FTCS-18)*, June 1988, pp. 15–22.
- [12] M. R. Lyu and A. Avižienis, "Assuring design diversity in n-version software: A design paradigm for n-version programming," *Dependable Computing and Fault-Tolerant Systems*, vol. 6, pp. 197–218, 1991.
- [13] D. E. Eckhardt and L. D. Lee, "A theoretical basis for the analysis of redundant software subject to coincident errors," National Aeronautics and Space Administration (NASA), Tech. Rep. NASA-TM-8636919850015006, 1985.
- [14] S. S. Brilliant, J. C. Knight, and N. G. Leveson, "The consistent comparison problem in n-version software," *IEEE Transactions on Software Engineering*, vol. 15, pp. 1481–1485, 1989.
- [15] D. Eckhardt and L. Lee, "A theoretical basis for the analysis of multiversion software subject to coincident errors," *IEEE Transactions on Software Engineering*, vol. 11, pp. 1511–1517, 1985.
- [16] B. Littlewood, P. P. L., and Strigini, "Modeling software design diversity: a review," *ACM Computing Surveys*, vol. 33, pp. 177–208, June 2001.
- [17] D. Partridge and W. Krzanowski, "Distinct failure diversity in multiversion software," Department of Computer Science, University of Exeter, U.K., Tech. Rep., 1997.
- [18] B. Littlewood and L. Strigini, "Redundancy and diversity in security," in *9th European Symposium on Research Computer Security, (ESORICS 2004) LNCS 3193*. Springer, 2004, pp. 423–438.
- [19] A. N. Bessani, R. R. Obelheiro, P. Sousa, and I. Gashi, "On the effects of diversity on intrusion tolerance," Department of Informatics, University of Lisbon, DI/FCUL TR 08–30, December 2008.
- [20] Y. Deswarte, K. Kanoun, and J.-C. Laprie, "Diversity against accidental and deliberate faults," in *Computer Security, Dependability, and Assurance: From Needs to Solutions*. IEEE Press, 1998, pp. 171–181.
- [21] M. K. Joseph and A. Avižienis, "A fault tolerance approach to computer viruses," in *Proceedings of the 1988 IEEE conference on Security and privacy*, ser. SP'88. Washington, DC, USA: IEEE Computer Society, 1988, pp. 52–58.
- [22] E. Barrantes and S. Forrest, "Increasing communications security through protocol parameter diversity," in *Proceedings of the XXXII Latin-American Conference on Informatics (CLEI 2006)*, August 2006.
- [23] A. D. Keromytis, "Randomized instruction sets and runtime environments past research and future directions," *IEEE Security and Privacy*, vol. 7, pp. 18–25, 2009.
- [24] J. Xu, Z. Kalbarczyk, and R. K. Iyer, "Transparent runtime randomization for security," in *In proceedings of the 22nd International Symposium on Reliable Distributed Systems (SRDS'03)*, 2003.
- [25] S. Bhatkar and R. Sekar, "Data space randomization," in *Detection of Intrusions and Malware, and Vulnerability Assessment, 5th International Conference, DIMVA 2008, Paris, France, July 10-11, 2008. Proceedings*, ser. Lecture Notes in Computer Science, vol. 5137. Springer, 2008, pp. 1–22.
- [26] M. Chew and D. Song, "Mitigating buffer overflows by operating system randomization," Carnegie Mellon University, Tech. Rep. CMU-CS-02-197, 2002.
- [27] A. Gherbi, R. Charpentier, and M. Couture, "Redundancy with diversity based software architectures for the detection and tolerance of cyber-attacks," DRDC Valcartier, Tech. Rep. TM 2010-287, 2010.
- [28] J. C. Knight and N. G. Leveson, "An experimental evaluation of the assumption of independence in multi-version programming," *IEEE Transactions on Software Engineering*, vol. 12, pp. 96–109, 1986.