

A Formal Framework for Positive and Negative Detection Schemes

Fernando Esponda, Stephanie Forrest, and Paul Helman

Abstract—In anomaly detection, the normal behavior of a process is characterized by a model, and deviations from the model are called anomalies. In behavior-based approaches to anomaly detection, the model of normal behavior is constructed from an observed sample of normally occurring patterns. Models of normal behavior can represent either the set of allowed patterns (positive detection) or the set of anomalous patterns (negative detection). A formal framework is given for analyzing the tradeoffs between positive and negative detection schemes in terms of the number of detectors needed to maximize coverage. For realistically sized problems, the universe of possible patterns is too large to represent exactly (in either the positive or negative scheme). Partial matching rules generalize the set of allowable (or unallowable) patterns, and the choice of matching rule affects the tradeoff between positive and negative detection. A new match rule is introduced, called *r-chunks*, and the generalizations induced by different partial matching rules are characterized in terms of the *crossover closure*. Permutations of the representation can be used to achieve more precise discrimination between normal and anomalous patterns. Quantitative results are given for the recognition ability of contiguous-bits matching together with permutations.

Index Terms—Anomaly detection, artificial immune systems, intrusion detection, negative detection.

NOMENCLATURE

l	Length of a string, a pattern or a packet (typically given in bits).
r	Size of the sliding window.
w	Window, a specification of r adjacent symbol positions in a string. When the window length r is understood, it suffices for w to specify only a start position.
t	Number of windows in a string of length l .
U	All possible strings of length l for a given alphabet.
S	Subset of U used to denote a sample of “normal” strings.
RS	Subset of U used to denote the complete set of “normal” strings.
$x(t)$	String or packet generated at time t .
G	Legitimate packet-generating process.

B	Malicious packet-generating process.
$y(t)$	Random variable indicating whether G or B generated $x(t)$.
\mathcal{H}	Hidden process governing the generation of packet $x(t)$.
sample(t)	Dynamic sample of G at instant t .
d	String detector.
Υ	Set of detectors.
$CC_r(S)$	Crossover closure generation rule of S for a given window size r .
$x[w]$	For $x \in U$, $x[w]$ is x projected onto window w .
dMx	d matches x under match rule M .
Scheme $_{ND}(\Upsilon)$	Set of strings x in U such that $(\forall \text{ windows } w)(\nexists d \in \Upsilon)(dMx)$.
Scheme $_{PD}(\Upsilon)$	Set of strings x in U such that $(\exists \text{ window } w)(\exists d \in \Upsilon)(dMx)$.
Scheme $_{PC}(\Upsilon)$	Set of strings x in U such that $(\forall \text{ windows } w)(\exists d \in \Upsilon)(dMx)$.
Scheme $_{NC}(\Upsilon)$	Set of strings x in U such that $(\exists \text{ window } w)(\nexists d \in \Upsilon)(dMx)$.

I. INTRODUCTION

IN anomaly detection, the normal behavior of a process is characterized by a model, and deviations from the model are called anomalies. In behavior-based approaches to anomaly detection, the model of normal behavior is constructed from an observed sample of normally occurring patterns. Models of normal behavior can represent either the set of allowed instances (positive detection) or the set of all anomalous instances (negative detection).

For most real problems, the set of positive instances is much smaller than the set of complementary ones. Thus, it would seem wise to derive a representation of the smaller set, rather than its enormous complement. However, there are several reasons why negative detection merits further consideration. First, it has been used successfully both in engineering applications and by naturally occurring biological systems. Second, if we assume a closed world, then from an information-theory perspective, the normal and abnormal sets both contain the same amount of information, which suggests that there might be equally compact representations of the negative patterns [13]. A third point about negative detection is that it allows the detection process to be distributed across multiple locations with virtually no communication required among the distributed components. This property allows several forms of distributed processing: checking small sections of a large object independently, independent detector sets (e.g., each one running on a separate machine), or

Manuscript received July 18, 2002; revised March 24, 2003. This work was supported by the National Science Foundation under Grants ANIR-9986555 and DBI-0309147, the Office of Naval Research under Grant N00014-99-1-0417, the Defense Advanced Projects Agency under Grant AGR F30602-00-2-0584, the Intel Corporation, the Santa Fe Institute, and the Consejo Nacional de Ciencia y Tecnología (México) under Grant 116691/131686. This paper was recommended by Associate Editor D. Cook.

The authors are with the Computer Science Department, University of New Mexico, Albuquerque, NM 87131 USA (e-mail: fesponda@cs.unm.edu).

Digital Object Identifier 10.1109/TSMCB.2003.817026

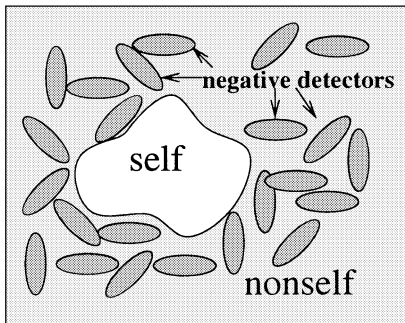


Fig. 1. Self/nonself discrimination. A universe of data points is partitioned into two sets: self and nonself. Negative detectors cover subsets of nonself.

evaluating each individual detector in a detector set independently. As the scale of anomaly-detection problems increases, the need for distributed processing is likely to grow. Indeed, the original inspiration for the negative-detection approach came from the natural immune system, which uses negative detection in a massively distributed environment—the body.

The negative-detection approach to anomaly detection, which was originally published as the negative-selection algorithm [1], [3], was modeled after a method used by the natural immune system to prevent autoimmunity [37]. In the immune system, certain cells known as T cells undergo a multistage maturation process in an isolated environment: an organ called the thymus. While in the thymus, T cells are censored against the normally occurring peptide patterns of the body called *self*. T cells that react with *self* are deleted in the thymus before they can become active and cause autoimmunity. The only T cells allowed to mature and leave the thymus are those that survive this censoring operation.¹ Such cells then circulate through the body freely and independently, eliminating any material that they can bind. Because of the censoring process, such material is implicitly assumed to be foreign and is known as *nonself*.

The translation of this mechanism into an algorithm for computers is straightforward. First, we assume that the anomaly-detection problem is posed as a set RS (real-self) of strings s , all of fixed-length l , of which we can access only a sample S at any given time. The universe of all l -length strings is referred to as U , and the set of anomalous patterns to be detected is the set $U - RS$. We consider all strings to be binary, although the analysis is generalizable to a larger alphabet. Candidate detectors are generated randomly and are censored against S ; those that fail to match any of the strings in S are retained as active detectors. Such detectors are known as *negative detectors*, and if S is a good sample of RS , each negative detector will cover (match) a subset of nonself. The idea is that by generating sufficient numbers of independent detectors, good coverage of the nonself set will be obtained. Fig. 1 shows the relationship of these sets pictorially.

This algorithm learns to distinguish a set of normally occurring patterns (*self*) from its complement (*nonself*) when only positive instances of the class are available. In the machine learning literature, this has been studied as the problem of

learning from examples of a single class or as concept learning from positive examples [25], [68]. Although in principle this learning problem is solvable for noise-free classes and certain formal domains, in many practical cases, the problem is known to be computationally intractable [27] or to lead to substantial overgeneralization [26]. The statistical community has examined the closely related problems of outlier detection and robust statistics, finding that the effectiveness of the learning system depends critically on domain-dependent assumptions about the distributions of self and nonself data. For example, in a computer security context, Lane employed time-series models of user behaviors and a prior bias against false alarms to differentiate intruders from authorized system users [28].

There are many well known change-detection and check-sum algorithms that solve a restricted form of the anomaly-detection problem, known as change detection. Here, *self* is known exactly, is small enough to be stored in a single location, remains constant over time, and can be unambiguously distinguished from *nonself*. However, for cases in which these assumptions do not hold, the discrimination task is more challenging, and in these situations, the negative-detection approach may be appropriate.

Since its introduction in [1], interest in negative detection has been growing, especially for applications in which noticing anomalous patterns is important, including computer virus detection [21], [22], intrusion detection [8], [15], [16], and industrial milling operations [4]. Recently, other categories of applications have been proposed, including color image classification [34] and collaborative filtering and evolutionary design [12]. Underlying all this work, however, is the question of negative versus positive detection: When is it advantageous to use a negative-detection scheme instead of the more straightforward positive-detection scheme?

To date, there has been little theoretical analysis of the relative merits of positive versus negative detection, and that is the primary goal of this paper. Such analysis is complicated by several factors, which we also address. Most important is the choice of representation and match rule. How is the set of normal patterns represented? How are anomaly detectors represented? What are the criteria for determining when a detector has detected a pattern? The analysis is further complicated if the *self* set is too large or too distributed to observe completely, if it is nonstationary (*self* sets that change over time), or if the detectors are changing over time.

In this paper, a formal framework is given for analyzing the tradeoffs between positive and negative detection schemes in terms of the number of detectors needed for maximum coverage. Earlier analyses focused on the cost of generating detectors (see Section II), rather than on the coverage provided by a complete set of detectors. We then discuss the need for partial matching rules and generalization, as well as what this implies for comparing positive and negative detection. We review the r -contiguous bits (rcb) matching rule, which is the most common match rule used with negative detection. We introduce a related match rule, called r -chunks, discuss its advantages, and characterize the generalizations that are induced by it in terms of a concept known as the *crossover closure* (see Fig. 3), giving the expected size of the crossover closure for r -chunks. We are

¹Mature T-cells have survived at least two other censoring operations: one involving genetic rearrangements and one involving positive selection.

then in a position to compare quantitatively the positive and negative detection schemes under r -chunks. Finally, we consider a method for increasing the diversity of coverage, known as permutation masks [8], [14], show that permutation masks reduce the number of nondetectable strings in the nonself set for any fixed detection threshold, and give quantitative results on how many permutations are required to maximize this effect.

One consequence of using an approximate scheme for representing a class of explicit instances (e.g., the self set) is that the class will implicitly contain some instances that were not in the original sample. If the approximation is a good one, then this is an advantage, and the approximation is said to *generalize* from the set of observed patterns to a useful class. In the case of anomaly detection, new observations that are similar to the original sample (self) are classified as part of self, and false positives are reduced. If the generalization is a poor one, new observations will frequently be misclassified, leading to high false-positives and/or false-negatives. In this paper, partial matching rules specify the approximation, and the r parameter controls how large the approximation is. We are interested in characterizing as precisely as possible the kinds of generalizations that are induced by our partial match rules. When we discuss the entire set covered by the approximation, we refer to it as the *generalization*. When we discuss individual strings in the generalization that were not in the original sample, we refer to them as *holes* [3]: strings not in S for which valid fixed- r detectors cannot be generated.

II. RELATED WORK

Outlier detection is a problem arising in many machine learning and data mining contexts. Abstractly, two families of approaches can be identified. Outlier detection with respect to probability distributions attempts to identify those events that are in the low-density regions of the probability distribution governing the generation of normal events [49], [50], [52], [58]. Under this approach, which we refer to as statistical anomaly detection, the degree of suspicion attached to an event is inversely proportional to the frequency with which the event has been observed historically.

A second family of approaches to outlier detection attempts to define a measure of distance on the event space [51], [54], [59], [60], [62], [64]. Common distance measures include the classic L -norm and Hamming distances, Manhattan distance (also known as rectilinear, which generalizes Hamming distance to nonbinary coordinate systems), and the vector cosine measure of Salton [65]. Under this approach, the degree of suspicion attached to an event is directly proportional to the distance of the event from, for example, the nearest observed normal event or the center of mass of clusters of normal events.

The main difficulty in applying statistical anomaly detection is that the probability distribution is not known exactly and must be estimated from a sample. Often, this sample is sparse, making density estimation highly error prone. In many applications, the majority of events we must judge do not appear in the sample at all and, hence, must be classified as “most suspect.” The main difficulty in applying distance criteria is in the construction of a measure that reflects a useful metric of similarity. A poor choice of measures results in meaningless classifications.

Often, either implicitly or explicitly, the approaches are combined [48], [61]. For example, data reduction transformations, such as feature selection and value aggregation, collapse distinct events into clusters of events that are to be treated as if they were indistinguishable [57], [66]. An event is assigned the composite density of the cluster into which it falls. In this way, many events are clustered together, all of which are classified as either suspect or not based on the cluster’s composite density. The most suspect events will cluster with no previously encountered events (forming a low density cluster of its own).

Statistical anomaly detection is often used on the problem of intrusion detection in computer security, for example, [15], [23], [35], [41]–[45], [55]. Intrusion-detection systems vary widely, but they all seek to protect an information system (e.g., a single computer, a database, a server, or a network of computers) from violations of the system’s security policy. In the case of anomaly intrusion-detection systems, protection is provided by building a model of the normal (legal) operation of the system and treating discrepancies from the model as anomalies. The model of normal behavior can be based on any observable behavior of the system. Audit logs, patterns of network traffic, user commands, and system calls are all common choices. Such an approach can work well if the anomalies in normal operation are well-correlated with security violations. The extent to which this is true is a topic of current debate in the intrusion-detection community.

In [1], the negative-selection algorithm was introduced as a method for representing and generating distributed sets of detectors to perform change detection. The negative-selection algorithm was demonstrated and analyzed in the context of the rcb match rule [38], [39], which is a partial-matching rule between two strings similar to Hamming distance. In that paper, a probabilistic expected-case analysis of the algorithm’s performance was given for the case where the protected data and the detectors are assumed to be static, and the protected data are randomly generated. The analysis relied on the assumption that independently generated detectors would have independent detection capabilities. Use of the algorithm was illustrated with the problem of computer virus detection, and it was noted that observed performance on real data sets was often significantly better than that predicted by the expected-case analysis. This improved performance was ascribed to regularities in the protected data (self).

References [3] and [13] reported a linear time algorithm for generating detectors (linear in the number of protected strings) that uses $O(2^r(l-r)^2)$ space, where r is a threshold value and l is the length of the strings. This new algorithm addressed an important limitation as the cost of generating detectors using the original algorithm increased exponentially with the size of the protected data. Limits to the precision of coverage that is possible under most plausible matching rules were also identified. These limits arise because partial matching rules obscure the boundary between the self and nonself sets. In [3], the number of such gaps, called holes, for certain partial matching rules were counted, and a lower bound on the size of the detector set needed to achieve a given probability of detecting abnormal strings was estimated. Following up on this work, Wierzchon developed a low space-complexity algorithm for generating an optimal

repertoire of detector strings [17]–[19] and included an analysis of the holes in coverage associated with his algorithm. In this work, it was assumed that both the protected data and the detectors are unchanging through time. Likewise, both projects focused on the problem of generating detectors efficiently, rather than the problem of how many detectors are needed for maximal performance. The latter question is the one addressed in this paper.

Lamont and Harmer applied negative selection to the problem of computer virus detection, focusing on different match rules [21]. They developed a prototype implementation to demonstrate that the approach is both scalable and effective. As part of their work, they evaluated 12 different match rules including Hamming distance and rcb, concluding that the Rogers and Tanimoto rule (an extension of Hamming distance) was the best fit for their application because it did the best job of balancing the generality and specificity of each individual detector.

Hofmeyr [7], [8], [14] described a network intrusion-detection system that incorporated a modified form of the negative-selection algorithm, together with several other mechanisms, including the use of permutation maps to achieve diversity of representation. His experiments were the first in this line of work to consider dynamically changing data sets (e.g., when the definition of normal behavior is either incomplete or nonstationary) and a distributed environment. The setting was network intrusion detection on a local area network, and the detectors monitored the flow of TCP SYN packets through the network. The original negative-selection algorithm was modified to accommodate dynamic data sets; detectors were generated asynchronously and allowed to undergo negative selection in the live environment of the network. Thus, the intrusion-detection system featured rolling sets of detectors in the sense that detectors generated at different times were potentially exposed to slightly differing patterns of self.

Kim and Bentley also implemented a network intrusion-detection system which incorporates negative selection [16]. Their implementation used a significantly more complicated representation than [14] with a much larger alphabet, shorter strings, and more general matching thresholds (lower r in r -contiguous bits matching). Like Hofmeyr, they used the original random-generation method of producing detectors. They reported difficulties generating valid detectors and concluded that negative selection is infeasible on its own for realistic network intrusion detection. Balthrop *et al.* [10] propose several explanations for their results, based largely on an analysis of the parameter settings that were selected. Williams *et al.* [15] report positive results when using negative selection on the network intrusion detection problem. Their results and approach more closely parallel that taken in Hofmeyr's work. Their implementation goes beyond [14] and [16], however, in that it monitors UDP and ICMP packets as well as TCP packets.

Dasgupta and Gonzalez used negative detection on the network intrusion-detection problem [33]. They worked with the Lincoln Labs data set [29], used a real-valued representation, and divided the Lincoln Labs data into overlapping intervals. The paper is interesting because it compares the performance of positive and negative characterization of self, reporting that the positive approach is more accurate but also computationally

more expensive. It is not obvious how meaningful the comparison is, however, because the positive and negative methods are quite different. The positive-detection method simply memorizes a set of training (attack-free) points and classifies any test point as “normal” if it is within a given distance of any memorized self-point. The negative-detection method uses a genetic algorithm to evolve complex rules for detecting nonself. Thus, the study could more properly be viewed as comparing evolved rules with a simple memory-based approach.

In addition to intrusion detection, several authors have been interested in using negative detection for other applications. Sathyanath and Sahin used negative detection for color image classification of finished wooden components: specifically, kitchen cabinets [34]. Bradley and Tyrrell describe a hardware implementation using field programmable gate arrays (FPGAs) [30]–[32]. Their detectors monitor transitions, together with some input/output information, in finite-state machines to detect faults. The implementation uses negative detection and rcb match rules. Detectors are generated offline, and the entire self set is known at the time detectors are being generated. Chao and Forrest [12] describe a negative-detection approach to interactive search algorithms, in which subjective evaluations drive the exploration of large parameter spaces. Their algorithm learns what parts of the search space are *not* useful, based on the negative-detection strategy of the natural immune system. The algorithm is capable of finding consensus solutions for parties with different selection criteria.

Sometimes, the set of anomalous patterns is known, or thought to be known, exactly, and a set of signatures is used to cover the dangerous parts of nonself precisely. Examples of this approach are signature-based scanners for intrusion detection [46] and commercial virus-scanners.

III. STATISTICAL MODELING FRAMEWORK

In this section, we cast the anomaly intrusion-detection problem in a statistical framework. Although the unrestricted version of the framework is quite complex, it is nevertheless useful to begin with the overall view. This helps us make rigorous certain notions which are required later in the paper, and it allows future work to evolve toward meeting broader objectives.

Our basic unit of analysis is a length l binary string, called a packet. At each time step t , a packet $x(t)$ is generated. We assume that $x(t)$ is generated either by the legitimate process G or the malicious process B . In reality, there could be any number of subprocesses contributing to each of these processes (e.g., many legitimate users with different behaviors), but we proceed as if there were a single G and a single B .

More precisely, a sequence $x(1), x(2), \dots, x(t)$ of length l binary strings (packets) is generated. A hidden process \mathcal{H} governs the generation of the packet $x(t)$ at time t . Random variable $y(t) = 0$ if process G generates $x(t)$ and $y(t) = 1$ if process B generates $x(t)$. In general, it may be the case that the distribution of random variable $y(t)$ depends on the sequences $x(1), x(2), \dots, x(t-1)$ and $y(1), y(2), \dots, y(t-1)$ of random variables, as well as on the identity of the particular time step t . In the remainder of this paper, however, we assume that the

distribution of $y(t)$ is independent of t and of $y(t')$ and $x(t')$ for $t' < t$. Thus, $Pr\{y(t) = 0\}$ and $Pr\{y(t) = 1\}$, which are also referred to as the prior probabilities of the processes $P(G)$ and $P(B)$, respectively, are unchanging over time.

There are several (typically unknown) next packet distributions of potential interest; estimating these distributions, under various assumptions to be introduced below, is a major focus of this paper. The pair of distributions

$$\begin{aligned} Pr\{x(t)|y(t) = 0, x(t-1), x(t-2), \dots, x(1) \\ y(t-1), y(t-2), \dots, y(1)\} \\ Pr\{x(t)|y(t) = 1, x(t-1), x(t-2), \dots, x(1) \\ y(t-1), y(t-2), \dots, y(1)\} \end{aligned}$$

specify next packet distributions under G and B , respectively, conditioned on knowledge of the first $t-1$ packets $x(t')$ generated, plus a specification (the values of $y(t')$) of which process generated each of these $x(t')$. Estimating these distributions from the data amounts to a supervised learning problem since the $y(t')$ label the training data $x(1), x(2), \dots, x(t-1)$, indicating which process generated each packet. Note that these probabilities are related via Bayes theorem to the posterior distribution on which process G or B generated the current packet $x(t)$. See Section III-B for details. Note also that the mixture distribution for the next packet generated is

$$\begin{aligned} Pr\{y(t) = 0\} \cdot Pr\{x(t)|y(t) = 0, x(t-1), x(t-2), \dots \\ x(1), y(t-1), y(t-2), \dots, y(1)\} \\ + Pr\{y(t) = 1\} \cdot Pr\{x(t)|y(t) = 1, x(t-1), x(t-2), \dots \\ x(1), y(t-1), y(t-2), \dots, y(1)\}. \end{aligned}$$

In the applications studied in this paper, we have no direct knowledge of the values of the random variables $y(t')$, that is, we do not know with certainty which of the processes G or B generated the previous packets in the sequence, leaving us with an unsupervised learning problem. In particular, we will be concerned with estimating the pair of distributions

$$\begin{aligned} Pr\{x(t)|y(t) = 0, \mathcal{H}, x(t-1), x(t-2), \dots, x(1)\} \\ Pr\{x(t)|y(t) = 1, \mathcal{H}, x(t-1), x(t-2), \dots, x(1)\} \end{aligned}$$

when we have knowledge of the process priors $P(G)$ and $P(B)$ (that is, knowledge of the distribution of the hidden process \mathcal{H}), or with estimating the pair of distributions

$$\begin{aligned} Pr\{x(t)|y(t) = 0, x(t-1), x(t-2), \dots, x(1)\} \\ Pr\{x(t)|y(t) = 1, x(t-1), x(t-2), \dots, x(1)\} \end{aligned}$$

when we do not have knowledge of the process priors.

For some problems, the ultimate objective is to identify short temporal groupings of anomalous packets or the time intervals most likely to contain packets generated by B . Many intrusion-detection problems have this quality, and the above formulation supports such an analysis. However, there are important special cases in which the analysis shifts from sequences of packets to single packets, and these are the focus here.

The next packet distribution of each process may depend on the current time step t , as well as on the pattern of packets generated at time steps prior to t . This allows the possibility that each of G and B is nonstationary, where its distribution depends both

on the actual time step t and the identity of packets previously generated by both processes. This nonstationarity might reflect, for example, that we have the following.

- Composite user behavior changes over time (different populations of users).
- A user's behavior at time t depends on behaviors at times $t' < t$ (a single user's action depends on his or her prior actions).
- The nature of attacks changes over time.
- There are changes in the environment (e.g., a new computer added to a network).

We abstract away the question of whether and how the true distributions of G and B change over time by specifying our posterior of the next packet $x(t)$ generated at any instant t in time, given some dynamic sample(t) of G at instant t . That is, in this paper, we focus on situations in which it suffices to replace

$$\begin{aligned} Pr\{x(t)|y(t) = 0, \mathcal{H}, x(t-1), x(t-2), \dots, x(1)\} \\ \text{or } Pr\{x(t)|y(t) = 0, x(t-1), x(t-2), \dots, x(1)\} \\ \text{by } Pr\{x(t)|y(t) = 0, \text{sample}(t), t\} \end{aligned}$$

where $\text{sample}(t)$ is an (unordered) set of packets. $\text{sample}(t)$ may include a subset of the packets $x(t')$, $t' < t$ restricted to packets deemed to have been generated by G (i.e., packets $x(t')$ for which it is deemed likely that $y(t') = 0$) and perhaps further restricted by other criteria (e.g., recency, random sampling). $\text{sample}(t)$ may also be seeded with an initial population of packets.

$\text{sample}(t)$ reflects our current state of knowledge of G at time t , and the next packet generated by G is dependent only on this knowledge—given $\text{sample}(t)$, $x(t)$'s distribution is conditionally independent of all other information (e.g., of other previously generated packets). As $\text{sample}(t)$ changes with time (see below), the next packet distribution changes. Whether or not process G actually is changing over time (or whether, for example, our state of knowledge simply is changing) is immaterial to the model. Further, because $\text{sample}(t)$ is unordered and assumed to be generated by or otherwise representative of G , the detection question of interest is to decide whether the next packet $x(t)$ is generated by G or B ; our model does not address the problem of deciding whether a subsequence embedded in $x(1), x(2), \dots, x(k)$ is generated by G or B .

In the remainder of this paper, we consider a simple family of distributions $Pr\{x(t)|y(t) = 0, \text{sample}(t), t\}$ as our modeling of how G generates its next packet. Intuitively, we assume a distance measure between an arbitrary packet p_k and the current $\text{sample}(t)$. The closer a packet p_k is to $\text{sample}(t)$, the greater the probability $Pr\{p_k|y(t) = 0, \text{sample}(t), t\}$. Many notions of distance are possible. For example, the distance of a packet p_k from $\text{sample}(t)$ could be p_k 's minimum Hamming distance to a member of $\text{sample}(t)$. Note that this measure does not distinguish between members of $\text{sample}(t)$ (assigning each such string a distance of 0) but does differentiate between packets not in $\text{sample}(t)$. Alternatively, distance could be based on frequency counts; the “distance” of p_k from $\text{sample}(t)$ could be the reciprocal of the number of instances of p_k in $\text{sample}(t)$, with $1/0$ defined to be some large constant C . This distance measure fails to distinguish between packets not in $\text{sample}(t)$ (assigning

each such string a distance of C) but does differentiate among packets in $\text{sample}(t)$. Of course, one could develop a hybrid of these two measures in which a distance based on Hamming distance is used for strings outside of $\text{sample}(t)$, and a distance based on frequency counts is used for strings in $\text{sample}(t)$.

We adopt a simple categorical division into “similar to $\text{sample}(t)$ ” versus “dissimilar from $\text{sample}(t)$ ” and distinguish these distance categories by means of a *generation rule* that attempts to characterize the underlying set from which $\text{sample}(t)$ is likely drawn. Unlike the two distance measures mentioned earlier, a generation rule allows only the distinction between likely and unlikely strings under G , which simplifies the detection task considerably. Experimental studies have shown that these simplifications often capture sufficient detail of process behavior to provide effective detection [5], [8].

Definition: A generation rule Q is a mapping from a set S of length l strings to a set $Q(S)$ of length l strings containing S . This is the generalization discussed in Section I.

Given a generation rule Q , we wish to specify a distribution to reflect that each member of $Q(\text{sample}(t))$ is more likely to be generated at time t by G than is each nonmember. Several simple alternatives exist, though care must be taken to ensure a proper probability distribution results. For example, for some applications, it may be appropriate to have the probability of the set $Q(\text{sample}(t))$ be some large multiplicative factor f of the probability of the complement of this set. We can construct a probability distribution capturing this relative frequency by first setting $p'_1 = f/|Q(\text{sample}(t))|$, $p'_2 = 1/|U - Q(\text{sample}(t))|$ and then (noting $p'_1 \cdot |Q(\text{sample}(t))| + p'_2 \cdot |U - Q(\text{sample}(t))| = f + 1$) normalizing to obtain proper point probabilities for individual strings, that is

$$\begin{aligned} Pr \{x(t)|y(t) = 0, \text{sample}(t), t\} \\ = \begin{cases} \frac{p'_1}{(f+1)}, & \text{if } x(t) \in Q(\text{sample}(t)) \\ \frac{p'_2}{(f+1)}, & \text{if } x(t) \notin Q(\text{sample}(t)). \end{cases} \end{aligned}$$

Assuming there are more strings in the complement of $Q(\text{sample}(t))$ than in the set itself, the point probability of each member string would be greater than that of each nonmember string. Alternatively, in other contexts (e.g., when $Q(\text{sample}(t))$ is larger than its complement), it would be necessary to specify directly that the common point probability of each member of $Q(\text{sample}(t))$ is some large multiple of the common point probability of each nonmember, and then again normalize to proper probabilities by taking into account the number of members and nonmembers.

Whatever method is used to set the distribution, we assume $Pr \{x(t)|y(t) = 0, \text{sample}(t), t\}$ is specified by

$$\begin{aligned} Pr \{x(t)|y(t) = 0, \text{sample}(t), t\} \\ = \begin{cases} p_1, & \text{if } x(t) \in Q(\text{sample}(t)) \\ p_2, & \text{if } x(t) \notin Q(\text{sample}(t)) \end{cases} \end{aligned}$$

where p_1 and p_2 satisfy $0 \leq p_2 < p_1 \leq 1.0$, and $p_1 + p_2 = 1.0$.

As a final comment, we note that under this modeling, although the generation rule Q and measure of distance (e.g., member or not of the set $Q(S)$) are assumed fixed and known exactly (as opposed to being dependent on the data), the probabilities assigned to $x(t)$ under G do change as the sample

changes with time. In particular, since $\text{sample}(t)$ typically is dependent on previously observed packets—evolving with our knowledge of the true image of process G —our distribution of the next packet generated changes as the data changes.

A. Example Generating Rules

We now present three sample generation rules that can be used to characterize the underlying set from which $\text{sample}(t)$ is drawn. The first is Hamming radius, which generalizes $\text{sample}(t)$ to strings within a given Hamming distance of strings in the sample. Next, we consider crossover closure, which generalizes $\text{sample}(t)$ to strings that can be constructed by pasting together the windows of strings appearing in $\text{sample}(t)$. Finally, we examine n -gram matching, which generalizes $\text{sample}(t)$ to strings whose n -grams have been observed in the $\text{sample}(t)$.

1) *Example: Hamming Radius:* The Hamming distance measure discussed above can be simplified to a categorical measure based on a threshold T . In particular, the generation rule $Hd_T(S)$ is defined to be the set of strings within a Hamming distance of T or less to some member of S . The use of this generation rule in detection is considered in Section VI. Hamming distance match rules have been used with negative detection in the context of immunological modeling, e.g., [6].

2) *Example: Crossover Closure:* The generation rule we consider in the most detail is called “crossover closure.” As we will show in subsequent sections, the crossover closure is closely related to the two match rules most commonly used in conjunction with negative detection.

Given a set S of strings and a fixed $1 \leq r \leq l$, the crossover closure $CC_r(S)$ of S is defined in terms of its length r windows (i.e., r consecutive string positions) as

$$CC_r(S) = \{u \in U | (\forall \text{ windows } w) (\exists s \in S) u[w] = s[w]\}.$$

In words, string $u \in U$ is in the crossover closure of S if and only if each of u 's t windows exactly matches the corresponding window of some member of S . When S is such that $CC_r(S) = S$, we say that S is *closed* under crossover closure. As we will see, the class of sets closed under crossover closure has a central role in our methods.

One interesting motivation for crossover closure comes from relational database theory [53], [63]. The strings of a sample $\text{sample}(t)$ can be viewed as the tuples of the current instance of a relation scheme $R(A_1, \dots, A_l)$, where each attribute A_i corresponds to packet position i and has domain $\{0, 1\}$. For example, the $\text{sample}(t) = \{0000, 1011\}$ can be represented as a relation scheme $R(A_1, A_2, A_3, A_4)$, whose current instance is shown as

$$\begin{array}{cccc} R(A_1, & A_2, & A_3, & A_4) \\ \hline 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1. \end{array}$$

For a variety of reasons (e.g., to reduce redundancy and enhance data integrity), it is often advantageous to represent a relation scheme as a decomposition into a collection of smaller relation schemes. Consider representing the scheme $R(A_1, A_2, \dots, A_l)$ as a decomposition into schemes $R_1(A_1, \dots, A_r)$, $R_2(A_2, \dots, A_{r+1}), \dots$, $R_i(A_i, \dots, A_{r+i-1}), \dots$, and $R_t(A_t, \dots, A_l)$. The instance

of each R_i is the projection of R onto R_i or, equivalently, is exactly the set of strings comprising the i^{th} length r window of $\text{sample}(t)$. In the previous example, taking $r = 2$, the instances of the R_i are as follows:

$$\begin{array}{cc|cc|cc} R_1(A_1, & A_2) & R_2(A_2, & A_3) & R_3(A_3, & A_4) \\ \hline 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 1 & 1. \end{array}$$

In order to reconstruct the original instance of R from these projections onto the R_i , one computes the natural join of the R_i . However, it is not always the case that the join of the projection recovers the original instance of R ; in fact, the join of the projected instances is precisely the crossover closure of the set of tuples (strings) in the original instance of R . In our example, the join $R_1|X|R_2|X|R_3$ of the instances shown above results in the following instance of R , which can be seen to be the crossover closure of the strings in the original instance of R .

$$R_1|X|R_2|X|R_3 = \begin{array}{cc|cc} R(A_1, & A_2, & A_3, & A_4) \\ \hline 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1. \end{array}$$

Relational database theory has exactly characterized when the projections of R join to recover the original instance (see Section VII-B and [53] and [63]). This is known as the lossless join condition. The lossless join condition thus also characterizes exactly when a set of packets is equal to its crossover closure. One interpretation of this correspondence is that since many naturally occurring collections of information do in fact satisfy the lossless join condition, it is plausible that in many contexts, the most likely set of strings to be generated by G is closed under crossover closure. When this is the case, $\text{sample}(t)$, if it is a representative sample of strings generated by G , can be interpreted as being a sample drawn from $CC(\text{sample}(t))$ or from a larger set containing $CC(\text{sample}(t))$ and itself closed under crossover. In such situations, it is appropriate to deem packets that are members of $CC(\text{sample}(t))$ as relatively likely under G .

3) *Example: n -Grams:* The n -gram matching rule has been used in a wide variety of settings including natural language processing [40], document classification [47], and program monitoring [5], [20], [24], [36] for intrusion detection. In the latter application, n -grams can be applied when the behavior of the protected process (for example, an executing computer program) can be represented as a sequence of letters taken from some alphabet. A sample of normal behavior S is a collection of such sequences, called a “trace” collected under normal operating conditions. Different traces can be of different lengths and are parsed using a sliding window of size n . Each substring of length n (the n -gram) is stored in the program’s profile.² Once the profile is complete, a new execution is analyzed in terms of the n -grams composing its trace. Any n -gram that occurs during execution, which does not appear in the profile, constitutes an anomaly.

²The n -gram method described here is a variant of the lookahead pairs method described in [2] and [9].

In order to characterize the generation rule for n -grams, we first restrict our attention to the generalization induced over traces of length l and create a directed acyclic graph (DAG) with $l - n + 1$ levels, where each level has a node for each distinct observed n -gram. A node with label P in level i is connected to a node with label Q in level $i + 1$ if the last $n - 1$ symbols of P match the first $n - 1$ symbols of Q (akin to the previous example). A similar graphical construction was used in [69] for determining multiple length n -grams. The set of protected traces of length l can be retrieved by traversing all possible paths, from the first to the last level, in this graph. We write the generalization induced for traces of size l under n -grams as $CC_n^l(S)$: the crossover closure at l (see Example 2). Finally, restricting traces to specific lengths, the generalization achieved by n -gram matching can be expressed as

$$CC_n^*(S) = \bigcup_{l \in L} CC_n^l(S)$$

where L is the set of acceptable lengths.

There are two differences between this generation rule and that presented in Section III-A2. In the latter, we considered all strings to be of the same specific length l ; this is a small assumption that could easily be relaxed. More importantly, in Section III-B2, we accept only those patterns that have been previously observed at a particular window position, whereas with n -grams, we disregard the information about window location.

B. Computing Posterior Probabilities on Intrusion Events

In general, deciding whether $x(t)$ is generated by G or B requires 1) distributions for B as well as G and 2) process priors $P(G)$ and $P(B)$. In previous work [55]–[57], we considered models in which $Pr\{x(t)|y(t) = 1, \text{sample}(t), t\}$ is uniform over all strings in U versus models for specific types of malicious behavior. Modeling $Pr\{x(t)|y(t) = 1, \text{sample}(t), t\}$ as uniform over all U is a conservative approach, and it is consistent with the view that we have little prior information about likely forms of attack. Modeling B as uniform places a constant c in the numerator of the likelihood ratio

$$\frac{Pr\{x(t)|y(t) = 1, \text{sample}(t), t\}}{Pr\{x(t)|y(t) = 0, \text{sample}(t), t\}}$$

and, hence, equates the posterior probability

$$\begin{aligned} Pr\{y(t) = 1|x(t), \text{sample}(t), t\} &= \frac{Pr\{x(t)|y(t) = 1, \text{sample}(t), t\} \cdot P(B)}{Pr\{x(t)|\text{sample}(t), t\}} \\ &= \frac{c \cdot P(B)}{c \cdot P(B) + Pr\{x(t)|y(t) = 0, \text{sample}(t), t\} \cdot P(G)} \end{aligned}$$

with the degree of $x(t)$ ’s rareness under G . (With the numerator of the expression for the posterior $Pr\{y(t) = 1|x(t), \text{sample}(t), t\}$ fixed at constant $c \cdot P(B)$, only $Pr\{x(t)|y(t) = 0, \text{sample}(t), t\}$ varies as $x(t)$ varies.) We identified this approach with what the literature commonly refers to as *anomaly detection* and reported experimental results with examples both when anomaly detection is sufficient and when it is not (and specific models of malicious behavior are necessary for satisfactory detector performance).

Here, we assume that B is uniform. With this assumption, we can either specify process priors $P(G)$ and $P(B)$ and a detection threshold on the posterior $Pr\{y(t) = 1|x(t), \text{sample}(t), t\}$, or, as in [56] and [57], we can simply prioritize the packets by the likelihood they were generated by B , that is, rank the packets by the values of their likelihood ratios; this gives an ordering equivalent to ranking by the posterior $Pr\{y(t) = 1|x(t), \text{sample}(t), t\}$, regardless of the values of the priors $P(G)$ and $P(B)$. With constant c in the numerator of the likelihood ratio, and the G process modeled as described in the previous section, ranking packets by their likelihood ratios reduces to partitioning packets into two classes: those that are most likely generated by G and those most likely not generated by G . We are interested in compact schemes that allow us to distinguish quickly between low- and high-probability anomalous packets. Under the family of distributions proposed above, this task is equivalent to distinguishing, for a given generation rule Q and current sample $\text{sample}(t)$, between packets contained in $Q(\text{sample}(t))$ and those outside this set. In the following sections, we present positive and negative detection schemes for addressing this problem. We also compare the time and space efficiencies of several variants of these schemes. All of these schemes use simple and efficient string matching rules for deciding whether or not a packet p_k is a member of $Q(\text{sample}(t))$.

IV. TAXONOMY OF DETECTION SCHEMES

This section explores several detection schemes based on the framework outlined above. Given a generation rule Q and a current sample $\text{sample}(t)$, an instance of a detection scheme (known simply as a detection algorithm) must be able to decide whether or not a packet p_k is a member of $Q(\text{sample}(t))$. That is, a detection algorithm can be viewed as protecting the subset $Q(\text{sample}(t)) \subset U$. Strings within $Q(\text{sample}(t))$ are regarded as likely to be part of self, whereas those in $U - Q(\text{sample}(t))$ are most likely to be anomalous. From the perspective of formal languages, a detection scheme is a language recognition model, and, analogous to the issues arising in the study of formal languages, we are interested in how to construct a detection scheme that can protect the class of languages implied by a particular generation rule of interest, as well as in determining the class of languages that various detection schemes can protect.

Many detection schemes can be constructed from the simple building blocks outlined so far—the simple string matching rules and the alternative interpretations associated with a match. Before proceeding to analyze specific cases of interest, e.g., [8], we first outline a systematic taxonomy for all the schemes. This taxonomy consists of three dimensions:

- the form of a single detector together with its match rule—a binary relation on detectors and packets, specifying when a single detector and single packet match;
- whether positive or negative detection is employed;
- disjunctive versus conjunctive matchings: Whether a single match is sufficient to trigger a decision (positive or negative) or whether multiple matches are required.

In the following sections, we study the relative power of the detection schemes that can be constructed from different choices

along the above dimensions. We first compare the classes of languages that can be protected by the schemes and then consider the resources (time and space) required by the schemes to protect their languages.

A. Form of Detector and Matching Rule

A single detector, as considered in this paper, is either

- 1) an element of U ;
- 2) an r -chunk, which is a length r binary string along with a specified window position.

When a detector d is an element of U , we are interested primarily in the rcb match rule, that is, for $d, x \in U$

$$dMx \leftrightarrow (\exists \text{ window } w)(d[w] = x[w]).$$

We refer to such detectors under the rcb match rule as rcb detectors. (In Section VI, we contrast the rcb match rule with Hamming distance, which is an alternative match rule for detectors in U .) If d is an r -chunk on window w , the matching rule we consider is, for $x \in U$

$$dMx \leftrightarrow d[w] = d.$$

B. Detector Operation

Let Υ be a set of rcb or r -chunk detectors, and assume the corresponding matching rule from the previous subsection. We now consider the interpretation given to a match between a packet and a detector. In particular, does a match mean the packet is “in the language” (positive detection, denoted by P) or “outside the language” (negative detection, denoted by N)? Further, do we require a single match anywhere in the packet (disjunctive matching, denoted by D), or a match in all packet window positions (conjunctive matching, denoted by C)? The four combinations of possible answers specify four corresponding detection schemes: PC , NC , PD , and ND .

Let α denote a choice of P or N and β a choice of C or D . With α and β specified, a collection Υ of detectors (either rcb or r -chunk) acts as a parameter, instantiating a detection scheme $\text{Scheme}_{\alpha,\beta}(\Upsilon)$. In particular, given the following interpretations for P , N , D , and C , $\text{Scheme}_{\alpha,\beta}(\Upsilon)$ exactly defines a set of allowable strings; $\text{Scheme}_{\alpha,\beta}(\Upsilon)$ protects, or recognizes, this set, flagging strings as anomalous that are outside $\text{Scheme}_{\alpha,\beta}(\Upsilon)$.

The language (subset of U) “accepted” by each of the four detection schemes when instantiated with a fixed set Υ of rcb or r -chunk detectors is defined as follows:

- 1) Negative disjunctive detection Scheme_{ND} : $\text{Scheme}_{ND}(\Upsilon)$ is the set of strings x in U such that $(\forall \text{ windows } w)(\nexists d \in \Upsilon)(dMx)$.
- 2) Positive disjunctive detection Scheme_{PD} : $\text{Scheme}_{PD}(\Upsilon)$ is the set of strings x in U such that $(\exists \text{ window } w)(\exists d \in \Upsilon)(dMx)$.
- 3) Positive conjunctive detection Scheme_{PC} : $\text{Scheme}_{PC}(\Upsilon)$ is the set of strings x in U such that $(\forall \text{ windows } w)(\exists d \in \Upsilon)(dMx)$.
- 4) Negative conjunctive detection Scheme_{NC} : $\text{Scheme}_{NC}(\Upsilon)$ is the set of strings x in U such that $(\exists \text{ window } w)(\exists d \in \Upsilon)(dMx)$.

C. r -Chunks Matching Subsumes rcb Matching

We say that $L \subset U$ is a language recognized by $\text{Scheme}_{\alpha,\beta}$ using r -chunks [rcb] detectors if there exists a set Υ of r -chunks [rcb] detectors such that $\text{Scheme}_{\alpha,\beta}(\Upsilon) = L$. From the perspective of language recognition, but not necessarily efficiency, r -chunk detectors subsume rcb detectors, as we now demonstrate.

Lemma 3.1: If d is an rcb detector, there exists a collection $T(d)$ of r -chunk detectors such that $\forall x \in U, d[w] = x[w] \leftrightarrow \exists t \in T(d), t = x[w]$.

Proof: Take $T(d)$ to be $\bigcup_{w_i} (d[w_i])$, and the result follows immediately. ■

To illustrate the construction of $T(d)$ used in the proof of the previous lemma, let $d = 1011$ be an r -contiguous bit detector with $l = 4, r = 2$. Then, $T(d) = \{d[1], d[2], d[3]\}$, as shown below, is the following equivalent set of r -chunk detectors:

$$\begin{array}{l} d: \quad \boxed{10111} \\ d[1]: \boxed{110} \\ d[2]: \quad \boxed{011} \\ d[3]: \quad \quad \boxed{111} \end{array}$$

It follows from the lemma that for each of the four detection schemes defined above, any set $L \subset U$ of strings that can be recognized with a set of rcb detectors can be recognized with some set of r -chunk detectors. In particular, we have the following theorem.

Theorem 3.1: If Υ is a set of rcb detectors, and α and β are any choices for detector dimensions as above, then $\text{Scheme}_{\alpha,\beta}(\Upsilon) = \text{Scheme}_{\alpha,\beta}(T(\Upsilon))$.

Proof: It follows from Lemma 3.1 that for any of the membership predicates $F_{\alpha,\beta}$ associated with any of the four detection schemes $\text{Scheme}_{\alpha,\beta}$ and for every string $x \in U$, $[\Upsilon \text{ satisfies } F_{\alpha,\beta} \text{ wrt } x] \leftrightarrow [T(\Upsilon) \text{ satisfies } F_{\alpha,\beta} \text{ wrt } x]$. ■

The converse of this theorem is not true in general, since r -chunk detectors have a finer granularity than do rcb detectors; a proper subset of $T(d)$ may match fewer strings in U than does the rcb d .

Example: Consider the detection scheme Scheme_{ND} , and let $l = 3$ and $r = 2$. Consider the pair of strings 011 and 010. We claim that whenever both strings 011 and 010 are in the language of $\text{Scheme}_{ND}(\Upsilon)$ for a set Υ of rcb detectors, then so too must be the string 110. To see this, note that an rcb detector matching 110 must either end in the pattern *10 or begin with the pattern 11*, but then, any way either of these patterns is completed (i.e., by specifying a bit for the *) results in an rcb detector matching one of 011 or 010. Consequently, if $\text{Scheme}_{ND}(\Upsilon)$ excludes 110 from the language it recognizes, it must exclude also at least one of 011 or 010. In contrast, if the r -chunks detector set Υ_{ch} consists of the single detector 11*, $\text{Scheme}_{ND}(\Upsilon_{ch})$ includes both 011 and 010 while excluding 110.

The results of the next section imply that the same result holds for Scheme_{PC} , that is, that the class of languages recognized by Scheme_{PC} using r -chunks detectors properly contains the class recognized by Scheme_{PC} using rcb detectors.

D. Detection Schemes and the Crossover Closure

The following theorem, which was established by simple set-theoretic arguments, helps clarify the relationships between

the classes of languages recognized by the various detection schemes.

Theorem 4.1: Let Υ be any set of rcb or r -chunks detectors, and let Υ' denote the complement of Υ relative to the universe of all rcb or r -chunks detectors over the same alphabet and of the same length as the detectors in Υ . Similarly, the complement of the set $\text{Scheme}_{\alpha,\beta}(\Upsilon)$ (which is a subset of U) is taken relative to U .

Then

$$\begin{aligned} \text{Scheme}_{ND}(\Upsilon) &= (\text{Scheme}_{PD}(\Upsilon))' \\ &= \text{Scheme}_{PC}(\Upsilon') = (\text{Scheme}_{NC}(\Upsilon'))' \\ &\subset \text{Scheme}_{ND}(\Upsilon')' = \text{Scheme}_{PD}(\Upsilon') \\ &= (\text{Scheme}_{PC}(\Upsilon))' = \text{Scheme}_{NC}(\Upsilon). \end{aligned}$$

Further, the subset containment is proper for some Υ .

It follows from $\text{Scheme}_{ND}(\Upsilon) = \text{Scheme}_{PC}(\Upsilon')$ that the class of languages recognized by Scheme_{ND} is identical to the class recognized by Scheme_{PC} , when either rcb or r -chunk detectors are considered. Similarly, it follows from $\text{Scheme}_{PD}(\Upsilon') = \text{Scheme}_{NC}(\Upsilon)$ that the class of languages recognized by Scheme_{PD} is identical to the class recognized by $\text{Scheme}_{NC}(\Upsilon)$, when either rcb or r -chunk detectors are considered.

As discussed in Section III, one of our most important detection applications is to protect the crossover closure of a sample $\text{sample}(t)$. There is a strong relationship between crossover closure and the schemes Scheme_{ND} and Scheme_{PC} when r -chunk detectors are used.

Let $\text{sample}(t)$ be any subset of U . Let W_P denote the set of windows present in $\text{sample}(t)$, that is, the union of the projections of $\text{sample}(t)$ onto each window (the window size is fixed and understood). Let W_N denote the set of windows not present in $\text{sample}(t)$, that is, $W_N = W_P'$. We then have the following theorem.

Theorem 4.2: $\text{Scheme}_{PC}(W_P) = \text{Scheme}_{ND}(W_P') = \text{Scheme}_{ND}(W_N) = CC(\text{sample}(t))$.

Proof: $x \in \text{Scheme}_{PC}(W_P) \leftrightarrow \forall_w x[w] \in W_P \leftrightarrow \forall_w \exists s \in \text{sample}(t), x[w] = s[w] \leftrightarrow x \in CC(\text{sample}(t))$. ■

Further, crossover closure exactly characterizes the class of languages recognized by Scheme_{PC} and Scheme_{ND} when r -chunk detectors are used. That is, we have the following theorem.

Theorem 4.3: The class of languages recognized by Scheme_{PC} and Scheme_{ND} when r -chunk detectors are used is exactly the class of sets closed under crossover closure.

Proof: If a set A is closed under crossover closure, i.e., $A = CC(A)$, then we can construct a set of detectors W_P from the windows present in A . It follows from Theorem 4.2 that $\text{Scheme}_{PD}(W_P) = \text{Scheme}_{ND}(W_P') = CC(A)$. By the definition of Scheme_{PC} , $x \in \text{Scheme}_{PC}(\Upsilon)$ if for all windows w $x[w] \in \Upsilon$, where Υ is an arbitrary set of detectors. We construct a set of detectors W_P by taking, for every window w and every $x \in \text{Scheme}_{PC}(\Upsilon)$, the projections $x[w]$ such that $x[w] \in \Upsilon \leftrightarrow x[w] \in W_P$. It follows that $x \in \text{Scheme}_{PC}(\Upsilon) \leftrightarrow x \in \text{Scheme}_{PC}(W_P)$. Finally, using Theorem 4.2, we have $\text{Scheme}_{PC}(\Upsilon) = \text{Scheme}_{PC}(W_P) = \text{Scheme}_{ND}(W_P') = CC(\text{Scheme}_{PC}(\Upsilon))$. ■

Results of the previous section imply that when rcb detectors are used, Scheme_{PC} and Scheme_{ND} can recognize only sets closed under crossover closure but not all sets closed under crossover closure since the class of languages recognized by these schemes when rcb detectors are used is properly contained by the class recognized when r -chunks are used.

Finally, consider the schemes Scheme_{PD} and Scheme_{NC} using r -chunks. These schemes do not recognize all sets closed under crossover closure but do recognize some sets not closed under crossover closure. To see this, we first observe that there exist sets S such that S is closed under crossover but S' is not. For example, let $l = 3$, $r = 2$, and $S = \{000\}$. $CC(S) = S$, and hence, S is closed under crossover. S' contains the other seven length 3 strings, and $CC(S') = U$.

Therefore, we have the following.

- 1) Scheme_{PD} and Scheme_{NC} cannot recognize all sets closed under crossover because if Scheme_{PD} recognizes a set S closed under crossover such that S' is not closed under crossover, Scheme_{ND}(Υ) = (Scheme_{PD}(Υ))' implies that Scheme_{ND} recognizes S' .
- 2) Scheme_{PD} and Scheme_{NC} can recognize some sets not closed under crossover because Scheme_{ND} can recognize a set S closed under crossover such that S' is not closed under crossover, and Scheme_{ND}(Υ) = (Scheme_{PD}(Υ)) l .

V. PARTIAL MATCHING AND GENERALIZATION

This section is devoted to Scheme_{PC} and Scheme_{ND} under the r -chunks matching rule. These schemes both recognize the sets closed under crossover, but they have different properties in terms of implementation requirements. These differences include distributivity, scalability, and algorithm efficiencies. This section explores the size of the detector set Υ and determines its expected size as a function of the number of strings in a randomly generated sample S .

A. Expected Number of Unique Detectors Under r -Chunks

Given a set S , it is straightforward to compute exactly how many detectors will be generated for maximal protection, both for the positive and negative detection schemes (Scheme_{PC} and Scheme_{ND}) using the obvious generation method. For Scheme_{PC}, it requires counting the number of distinct patterns for each of the t windows that comprise the strings in S , whereas for Scheme_{ND}, enumerating the distinct patterns that are not present in each window will result in the number of detectors. To provide an estimate of the average number of detectors, for both cases, as a function of the size of S , we note the following.

- The number of strings in U with a specific pattern in any given window of size r is 2^{l-r} .
- The probability of selecting at random one such string is $2^{l-r}/2^l = 2^{-r}$.
- Assuming r is small compared with l , we approximate the probability that a randomly generated self set of $|S|$ unique strings will contain a specific pattern by $[1 - \binom{|S|}{0}(2^{-r})^0(1 - 2^{-r})^{|S|} = 1 - (1 - 2^{-r})^{|S|}]$. The

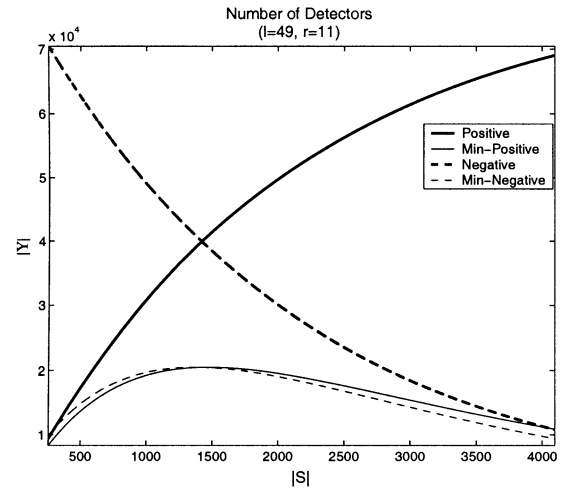


Fig. 2. Number of positive and negative detectors as a function of the size of the self set, assuming the self set was generated randomly. The plot shows both complete and reduced detector sets.

equation is approximate because it considers trials to be independent and sampling to take place with replacement.

- Following the previous item, the expected number of distinct patterns E_r , for a given window of size r , is approximated by $E_r \approx 2^r - 2^r(1 - 2^{-r})^{|S|}$.

The following equation denotes the expected size of a set of detectors having the property that each member matches one window and that all windows are matched:

$$E_{\text{pos}} = tE_r. \quad (1)$$

Conversely, the expected number of detectors possible for the negative detection scheme (i.e. detectors that do not match any window pattern in S) is given by

$$E_{\text{neg}} = t(2^r - E_r). \quad (2)$$

To establish when one scheme requires a smaller set of detectors than the other for maximal coverage, we determine the number of strings in S for which both schemes yield the same number of detectors, i.e., when $E_{\text{pos}} = E_{\text{neg}}$ (see Fig. 2):

$$E_{\text{pos}} = E_{\text{neg}} = t \left(2^r(1 - 2^{-r})^{|S|} \right) = t \left(2^r - 2^r(1 - 2^{-r})^{|S|} \right).$$

Solving for $|S|$

$$|S| = \frac{-\ln(2)}{\ln(1 - 2^{-r})} \approx (0.693)2^r. \quad (3)$$

Although we are primarily concerned with binary strings, it is worthwhile to note that for an arbitrary size alphabet \mathbb{A} , the above equation generalizes to $|S| \approx (0.693)|\mathbb{A}|^r$, which affects the point at which one scheme is preferable to the other, in terms of the size of the detector set.

In a worst-case scenario, Scheme_{PC} may require $t2^r$ detectors when $|S| \geq 2^r$. Similarly, Scheme_{ND} can also yield up to $t2^r$ detectors but only when there are no self strings whatsoever.

1) *Reduced Detector Set for Negative Detection:* The number of detectors needed in Scheme_{ND}, to protect $CC(S)$, can actually be smaller than the full set described above since significant redundancy amongst detectors may be present. We

examine this property in detail for the case of a binary alphabet. Number the t windows from left to right.

- Regardless of the composition of S , a detector must be generated for each pattern in the first window that is not present in S .
- For every window of size r , starting from the second window, and for every pair of patterns $w1 = v_i \dots v_{r+i-2}a$, $w2 = v_i \dots v_{r+i-2}\bar{a}$ in such a window, we have the following.
 - If both $w1$ and $w2$ are present in self, then we cannot generate any detector for them.
 - If neither is present, then we need not generate detectors for them since the preceding window will be missing its prefix, i.e., $bv_i \dots v_{r+i-2}$ or $\bar{b}v_i \dots v_{r+i-2}$, and a detector in the previous window will also match strings with such a pattern.
 - If only one is present, say $v_i \dots v_{r+i-2}a$, then we must generate detector $v_i \dots v_{r+i-2}\bar{a}$ since no string in S contains it.

With this in mind, the average number of detectors needed in the minimal set is given by

$$E_{minN} = 2^r - E_r + (l - r)(E_r - 2(E_r - E_{r-1})). \quad (4)$$

Following the previous rationale, we can also set an upper bound on the number of detectors required in the minimal set. The maximum number of self strings we can have without creating crossovers, thereby reducing the number of required detectors, will exhibit only one of every pair of patterns $w1 = v_i \dots v_{r+i-2}a$, $w2 = v_i \dots v_{r+i-2}\bar{a}$ for each window. This results from a maximum of 2^{r-1} distinct self strings for a detector set size of $t2^{r-1}$. Fig. 2 shows the plots of the expected number of detectors for both the full and reduced detector sets.

2) *Reduced Detector Set Size for Positive Detection:* Using similar reasoning to that of Section V-A1, we can find a significant amount of redundancy amongst detectors in the form of implicit matches. Consider the case where window $i + 1$ contains patterns $w1 = v_i \dots v_{r+i-2}a$, $w2 = v_i \dots v_{r+i-2}\bar{a}$; then, it must be that window i has either or both of the patterns ending with bits $v_i \dots v_{r+i-2}$, and therefore, a string matched in window i by one of these patterns will also be matched in window $i + 1$ by either $w1$ or $w2$. We call such a match an *implicit match* and eliminate $w1$ and $w2$ from the detector repertoire. The number of detectors in the resulting set can be expressed as

$$E_{minP} = E_r + (l - r)(E_r - 2(E_r - E_{r-1})). \quad (5)$$

The size of the sample S for which E_{minN} and E_{minP} yield the same number of detectors is the same as with the full repertoire, i.e., E_{neg} and E_{pos} (3).

One subtlety about this analysis is that if not all positive detectors are represented explicitly, then some additional information is required to identify implicit matches. This could be stored explicitly, requiring at the very least one bit per implicit match, or it could be derived at runtime by determining, once a match (or implicit match) at window i has been established, if both $w1$ and $w2$ are absent in window $i + 1$.³ A plot of E_{minP} for different

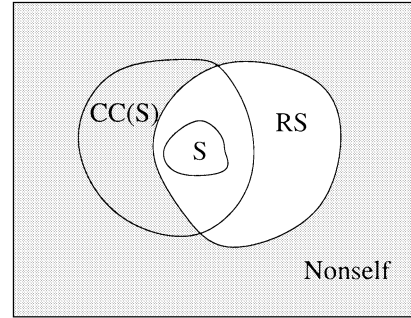


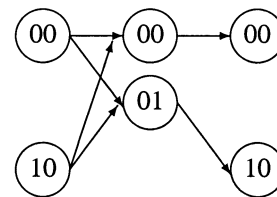
Fig. 3. Crossover closure given a sample S of real self data RS .

sample sizes is presented in Fig. 2 without regard to the extra information required to determine implicit matches.

B. Crossover Closure and Its Expected Size

One useful way to visualize the generalization induced by the r -chunks detectors is to construct DAG $G = (V, E)$, where $|V|$ is the number of distinct bit patterns of length r over all windows t of strings in S . Vertices are labeled by a pair (i, j) , where i represents the window number, and j stands for the bit pattern. $|E|$ is the number of edges, and edge $e((i, j), (i', j')) \in E \leftrightarrow i' = i + 1$ and the last $r - 1$ bits of pattern j match the first $r - 1$ bits of pattern j' , where $(i, j), (i', j') \in V$. Under this representation, the crossover closure is exactly the set of strings formed by all traversals of the graph from levels 1 to t .

Take, for instance, a self set S comprised of the following two strings $S = \{0000, 1010\}$ with $l = 4$, $r = 2$ (see Fig. 3). The corresponding graph $G = (V, E)$ is $V = \{(1, 00), (1, 10), (2, 00), (2, 01), (3, 00), (3, 10)\}$ and $E = \{((1, 00), (2, 00)), ((1, 00), (2, 01)), ((1, 10), (2, 01)), ((1, 10), (2, 00)), ((2, 00), (3, 00)), ((2, 01), (3, 10))\}$. A visual depiction of the graph (omitting the window number from the vertex labels) is⁴



Recovering the strings by traversing the preceding graph, we find that $CC(S) = \{0000, 0010, 1000, 1010\}$. In order to determine the size of $CC(S)$ for a randomly generated sample, we note that the number of substrings of length l that contain pattern $v_i \dots v_{r+i-1}$ in window w_i is double the number of substrings of length $l - 1$ that contain the pattern in the same window if there are strings in S that exhibit both $v_{i+1} \dots v_{r+i-1}a$ and $v_{i+1} \dots v_{r+i-1}\bar{a}$ in window w_{i+1} and stays the same if exactly one of these is present. In terms of our DAG representation, the number of paths that include a node with a given label doubles

³If the information is stored explicitly, the extra bits can be “recovered” as a decrease in execution time.

⁴The labels for window numbers could be recovered from the level by using a topological sort.

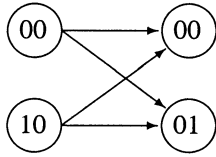
when such a node has two outgoing edges. We can write this as a recurrence on the number of windows t :

$$CC(t) = \begin{cases} E_r, & \text{if } t = 1 \\ 2CC(t-1)P(2) \\ + CC(t-1)(1-P(2)), & \text{otherwise.} \end{cases}$$

Solving the recurrence yields:

$$CC(t) = E_r (1 + P(2))^{(t-1)} \quad (6)$$

where $P(2)$ is the probability of a node having two outgoing edges. In order establish the value for $P(2)$, consider the following gadget:



The probability for a given edge to be present is approximated by $(E_{r+1}/2^{r+1}) = 1 - (1 - 2^{-(r+1)})^{|S|}$ and its absence by $((2^{r+1} - E_{r+1})/2^{r+1}) = (1 - 2^{-(r+1)})^{|S|}$.

Given that the likelihood of a node having only one outgoing edge is not independent of the probabilities related to the second node in the graph (a node at the same level differing only in the first bit position), we consider the probability $P^*(1)$ of either node having one outgoing edge:

$$P^*(1) = 4 \frac{E_{r+1}}{2^{r+1}} \left(\frac{2^{r+1} - E_{r+1}}{2^{r+1}} \right)^2.$$

Similarly, for a node in the gadget to have no outgoing edges (or for the node to be absent)

$$P^*(0) = \left(\frac{2^{r+1} - E_{r+1}}{2^{r+1}} \right)^4.$$

Finally, the probability of an individual node having two outgoing edges is given by

$$P(2) = 1 - \frac{1}{2} (P^*(0) + P^*(1)).$$

The number of holes $|H|$ can be derived by simply subtracting $|S|$ from $CC(t)$:

$$|H| = CC(t) - |S|. \quad (7)$$

The size of the generalization $CC(S)$ can be calculated from the size of S or from the size of the detector set (obtaining an estimate for $|S|$ from either (2) or (1) and substituting in (6)). It is important to note that the actual size will depend on the structure of the specific self set. Nevertheless, the analysis provides insight into its behavior (see Fig. 4) and enables us to ascertain the impact of allowing novel strings into the sample. This can be useful for determining, in a dynamic scenario, when (or at what rate) should detectors be added or deleted from the working set.

VI. PERMUTATIONS AND DIVERSITY

In [8] and [14], an additional mechanism was introduced to improve discrimination between self and nonself. This mechanism is loosely modeled after the diversity of major histocom-

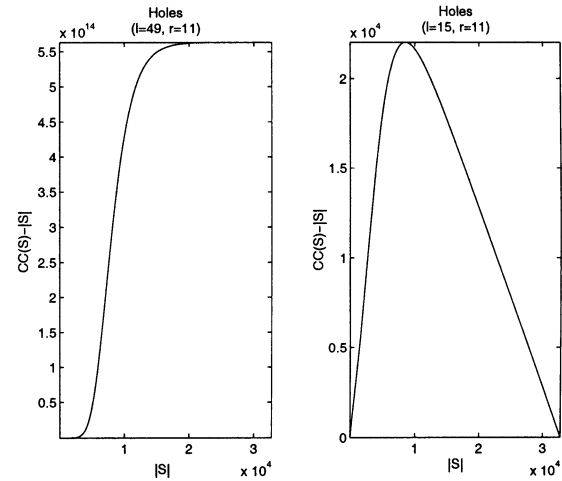


Fig. 4. Number of holes as a function of self-sample size for two parameter settings. The graph on the left illustrates how fast the generalization increases when r is small compared with the string length. The graph on the right shows how, after the generalization has reached its maximum, the number of holes slowly decreases as more strings are added to the sample S .

patibility complex (MHC) molecules used in the natural immune system. In the artificial setting, multiple permutations of the l -length strings (both data and detectors) are stored. This, combined with the contiguous-bits match rules, provides the system with *diversity* of representation and improves the discrimination abilities of the system.

A. Permutation Maps (π) for Reducing the Number of Holes

If H , which is the set of holes, is too large, accurate discrimination between RS and $U - RS$ will be unachievable. Hofmeyr introduced the concept of *permutation masks* as a way to control the size of H . A permutation mask is a reversible mapping that specifies a reordering of bits for all strings in U . In his network intrusion-detection data set, he reported an improvement in detection ability by about a factor of three, when permutation masks were used.⁵ By controlling the number of permutations used, he was able to achieve greater control over the r -contiguous bits generalization. One interesting property of this transformation is that if we consider all nonredundant permutation masks, the r -contiguous bits matching rule is able to protect strictly more sets than Hamming distance.

Theorem 1.1: The class of languages recognized by r -contiguous bits, with permutation masks, properly contains that recognized by Hamming distance.

Proof:

- Let x, d be strings of length l , where d is a detector.
- Let $\text{Rcb}(x, d) : \text{Strings} \times \text{Strings} \rightarrow \mathbb{N}$ be a function that returns the length of the longest run of contiguous bits that match between strings x and d . String d is said to match string x if $\text{Rcb}(x, d) \geq r$.
- Let $\text{Hd}(x, d) : \text{Strings} \times \text{Strings} \rightarrow \mathbb{N}$ be a function that returns the Hamming distance between strings x and d . String x is matched by a string d using the Hamming distance rule if d and x differ in, at most, $l - r$ bit positions, $\text{Hd}(x, d) \leq l - r$.

⁵Hofmeyr's dataset was collected using a variant of pure permutations that was computationally more efficient.

- Let $\pi(x)$ denote some permutation of string x .
- Let $L(\text{Rcb})$ denote the class of languages recognized by Rcb with permutation masks for a given r and a detector set Υ .
- $L(Hd)$ denote the class of languages recognized by Hd given r and a set of detectors Υ .

The class of languages recognized by Rcb, with permutation masks, contains the class of languages recognized by Hd :

- 1) $Hd(x, d) > l - r \rightarrow \text{Rcb}(\pi(x), \pi(d)) < r$ for all permutations π , and hence, $[d$ is a valid detector under $Hd]$ implies $[d$ matches no member of $\pi(L(\text{Rcb}))$ under Rcb and any permutation $\pi]$.
- 2) Further, if d matches a string x under Hd , we show there exists a permutation π^* such that d matches $\pi^*(x)$ under Rcb, that is, $L(Hd) \subset L(\text{Rcb})$. Label the bits that differ between x and d as $c_1c_2 \dots c_n$, $1 \leq n \leq Hd(x, d)$ and the bits that match as $c_{n+1}c_{n+2} \dots c_m$, $n + 1 \leq m \leq l$. We define π^* as $c_1c_2 \dots c_n c_{n+1} \dots c_l$. Clearly, π^* is a valid permutation, and $\text{Rcb}(\pi^*(x), \pi^*(d)) \geq r$. By (1) d is a valid detector for $\pi^*(L(\text{Rcb}))$ that matches $\pi^*(x)$. Hence, if L is a language recognized by some set of detectors under Hd , we have shown how to construct an associated set of permutations such that this same set of detectors recognizes exactly L also under Rcb. ■

In contrast, Hd does not recognize all languages recognized by Rcb with permutation masks because a detector d under Hd might match more than it matches under a single permutation and Rcb and thus fail to be a valid detector under Hd .

Example: Let $S = \{011110100\}$, $l = 3$, $r = 2$, and $h = 000$. Under Rcb, we can generate a valid detector $d = 001$ to match h , whereas under Hd , the potential detectors for h are 000, 001, 010, and 100, all of which have a Hamming distance less than two from h but also from all strings in S . Therefore, h is undetectable under Hd but detectable under Rcb with permutation masks; hence $L(\text{Rcb}) \not\subset L(Hd)$.

For those sets that Hamming distance can recognize, the Rcb rule may require multiple sets of detectors Υ , each with its own transformation π . This mechanism is suitable for a detection scheme that requires precise control over the generalization or that must be distributed.

To illustrate how the permutation masks can in fact reduce the number of holes, consider the self strings $s1 = 101$ and $s2 = 000$ with $l = 3$ and $r = 2$. There are two crossover strings $h1$ and $h2$ other than $s1$ and $s2$, for which no detector can be generated since neither string has a distinctive window that separates it from self:

$$\begin{aligned} s1 &= \boxed{101} & h1 &= \boxed{1101} \\ s2 &= \boxed{010} & h2 &= \boxed{0011} \end{aligned}$$

If we apply the permutation by which the third bit position is next to the first bit, as the following picture depicts, then we can generate two detectors $d1$ and $d2$ that will be able to match $h1$ and $h2$, respectively.

$$\begin{aligned} \pi(s1) &= \boxed{110} & \pi(h1) &= \boxed{1101} & d1 &= \boxed{101} \\ \pi(s2) &= \boxed{010} & \pi(h2) &= \boxed{0110} & d2 &= \boxed{011} \end{aligned}$$

TABLE I
UNDETECTABLE STRINGS UNDER EVERY PERMUTATION, FOR A GIVEN SET S
AND THE rcb MATCHING RULE

Self	010	001	100	100	001	010
	011	011	101	110	101	110
	100	100	010	001	010	001
	101	110	011	011	110	101
h	110	101	110	101	011	011
detector	11*	10*	11*	10*	01*	01*
templates	*10	*01	*10	*01	*11	*11

TABLE II
UNDETECTABLE STRINGS UNDER EVERY PERMUTATION, FOR A GIVEN SET
 S AND THE r -CHUNKS MATCHING RULE

Self	111	111	111	111	111	111
	010	100	010	100	001	001
	100	010	001	001	100	010
h	110	110	011	101	101	011
detector	11*	11*	01*	10*	10*	01*
templates	*10	*10	*11	*01	*01	*11

B. Completeness

We can now ask whether there are languages that cannot be recognized by rcb, even if permutation masks are used. That is, are there still strings in $\text{nonself}(U - RS)$ that go undetected even if every permutation mask is employed? The answer to this question is yes; as pointed out in [3], all practical match rules with a constant matching probability will exhibit holes.

We show, by construction, that there are languages that cannot be recognized under rcb matching with permutations: Let $l = 3$, $r = 2$ and $S = \{010011100101\}$, and consider the hole $h = 110$. Table I lists all the possible permutations of S and the substrings out of which a detector could be constructed for h under each permutation.

As can be seen from Table I, under each permutation, there are no available templates (that is, templates that do not match some string in S), and thus, neither rcb Scheme_{ND} nor rcb Scheme_{PC} can accept the strings in S while excluding h . Although we can construct such examples, it is difficult to characterize the set of all such holes under permutations.

Lifting the restriction that detectors be of length l enables a r -chunk detector to match the string h from the previous example. Nevertheless, there are languages that cannot be recognized under the r -chunks matching rule (Scheme_{PC} or Scheme_{ND}) augmented with permutation masks. Consider the following scenario: Let $l = 3$, $r = 2$, $S = \{111010100\}$, and $h = 110$. Table II lists all six permutations of S and h and the possible r -chunks detectors needed to match h .

In general, a string h cannot be detected by Scheme_{PC} or Scheme_{ND} under any permutation if all of its template combinations are in self. The template combinations of a string are the templates (strings with unspecified bits) that result from all the possible ways of specifying only r bits. For instance, the template combinations of $h = 110$ are $\{11^*, 1^*0, ^*10\}$, where * denotes “don’t care.” It is easy to see that only if one such template does not exist in S can there exist a permutation that places these bit positions contiguously, without creating the same contiguous bit pattern in S .

VII. DISCUSSION

The preceding sections developed a formal framework for studying tradeoffs between negative and positive detection and presented several theoretical results that we believe are important to a wide variety of practical problems. In this section, we explore some of the implications and potential extensions of our theoretical results.

A. Computer Security Applications

Over the past several years, a number of anomaly intrusion-detection systems have been developed for computer security, which explore different instances of the detection schemes described in this paper. These systems were reviewed in Section II; some use positive detection [2], [5], [9], and some use negative detection [1], [8]. Likewise, some of these systems use r -contiguous bits matching [1], [8], some use r -chunks [11], some use n -gram matching [5], and some use a variant of n -gram matching known as “lookahead pairs” [2], [9]. Although the negative-selection strategy of the immune system has received a great deal of attention, it should be noted that the immune system also uses positive selection, combining it advantageously with negative selection.

In some cases, we have good experimental evidence for preferring one scheme over another. For example, n -grams outperformed lookahead pairs on our early data sets in terms of absolute discrimination ability (unpublished). However, when efficiency matters, as in the case of online detection, then the lookahead-pairs method is a clear winner, paying a small penalty in terms of discrimination ability. A second example occurred when we compared rcb detection to r -chunks on a network intrusion-detection task [10], [11]. Such results are anecdotal in the sense that they were obtained by experimentally testing one or two methods on a limited number of data sets. Without some theory to guide us, it is difficult to determine how much a given result depends on the particular data set used and how much it depends on the choice of method. Likewise, it has been difficult to determine which aspects of a given detection scheme are most responsible for its success (or failure), e.g., negative detection, match rule, or parameter settings.

The results presented in this paper allow us to begin approaching such questions from a theoretical perspective. In particular, the notion of a crossover closure, and its relation to rcb and r -chunks matching, will allow us to understand more deeply when and why these matching methods are preferable to more familiar methods such as Hamming distance. Likewise, the closed-form expressions for detector set sizes, both for positive and negative detection, allow us for the first time to predict how large a problem must be before it pays to use a negative-detection scheme. Until now, the negative-detection approach has been somewhat of a curiosity, notorious as much for its immunological metaphor as for its demonstrated advantages over other methods. Results such as those presented here will form the basis of a more objective evaluation.

B. Relational Databases

Section III introduced the relationship between crossover closure and relational database theory. Specifically, we showed that

the crossover closure of a set S of strings is equivalent to the natural join of those strings (interpreted as tuples) projected onto a relational decomposition scheme. There are several aspects of this equivalence, which we find interesting.

First, we can characterize when a set of strings is closed under crossover closure in terms of well-studied structural properties known as database dependencies [53], [63]. A set of strings is closed under crossover closure exactly when it decomposes losslessly into the relation schemes

$$R_1(A_1, \dots, A_r), R_2(A_2, \dots, A_{r+1}), \dots, \\ R_i(A_i, \dots, A_{r+i-1}), \dots, R_t(A_t, \dots, A_L)$$

induced by its t r -length windows, as we showed in Section III-A. Although the structural characterization of a lossless join decomposition is complex in the general case, the specialized form of the above decomposition simplifies the condition considerably; the characterization reduces to the adherence to a collection of simple multivalued dependencies on the original relation scheme $R(A_1, A_2, \dots, A_L)$. In order for the decomposition to be lossless, one of the following multivalued dependencies on R must hold for each $i = 1, 2, \dots, (t-1)$:

$$A_{i+1}, \dots, A_{r+i-1} \twoheadrightarrow A_i$$

or

$$A_{i+1}, \dots, A_{r+i-1} \twoheadrightarrow A_{r+i}.$$

Note that, for example, the multivalued dependency $A_{i+1}, \dots, A_{r+i-1} \twoheadrightarrow A_i$ asserts that if R contains the tuples

$$a_1 \dots a_{i-1} a_i a_{i+1} a_{i+2} \dots a_{r+i-1} a_{r+i} \dots a_L$$

and

$$b_1 \dots b_{i-1} b_i a_{i+1} a_{i+2} \dots a_{r+i-1} b_{r+i} \dots b_L$$

then it must also contain the tuples

$$a_1 \dots a_{i-1} b_i a_{i+1} a_{i+2} \dots a_{r+i-1} a_{r+i} \dots a_L$$

and

$$b_1 \dots b_{i-1} a_i a_{i+1} a_{i+2} \dots a_{r+i-1} b_{r+i} \dots b_L.$$

This condition describes the semantics of collections of information that occur naturally within a vast number of diverse application domains. Whether such dependencies arise naturally among sets of strings we wish to protect remains to be investigated.

Taking the connection one step further, we can investigate the correspondence between the schemes Scheme_{ND} and Scheme_{PC} , which recognize the class of sets closed under crossover closure and query systems for relational databases. In particular, if one is able to recognize with a set of detectors (possibly significantly pruned to remove redundancy as described in Section V-A2) a given set closed under crossover closure, one can also represent by the same means any instance of a relation scheme obeying the dependencies specified above. The negative detection scheme Scheme_{ND} would be most appropriate for membership queries, whereas Scheme_{PC} would be appropriate for queries requiring enumeration of

the member tuples. Conversely, known results from database theory regarding minimal representations translate to lower bounds on time and space requirements for these detection schemes.

C. Dynamic Samples

As mentioned in Section III, we expect our samples—and hence our distribution of the next packet generated by G —to change with time. Some factors contributing to this include the following.

- 1) If $x(t)$ is not in $Q(\text{sample}(t))$, it may become part of $\text{sample}(t+1)$, and $Q(\text{sample}(t+1))$ then properly contains $Q(\text{sample}(t))$, and consequently, the distribution of the string generated at time $t+1$ differs from that at t . This assumes we have a mechanism for deciding that $x(t)$ is generated by G , despite the fact that we flag it as a low probability string (e.g., it passes some investigative step).
- 2) We might want to delete packets from the sample that have not occurred recently. This might or might not alter $Q(\text{sample}(t))$ and, hence, the distribution of the string generated at time $t+1$.

Extending our analytical treatment to include dynamically changing samples raises a number of questions. In principle, when a new sample is created at time t , a new set of detectors could be generated *de novo*. However, this is likely to be extremely inefficient. An area of future investigation is to analyze the techniques used to update the detector set in response to the specific types of perturbations in the sample.

D. Distance Measures

If we can define a metric over a match rule, then our detection schemes can be rewritten in terms of a distance measure. For instance, using r -chunks

- $\text{Scheme}_{PD}(W_P) = \{x | d(x, S) < t\}$;
- $\text{Scheme}_{PC}(W_P) = \{x | d(x, S) = 0\}$;
- $\text{Scheme}_{ND}(W_N) = \{x | d(x, S) = 0\}$;
- $\text{Scheme}_{NC}(W_N) = \{x | d(x, S) < t\}$;

where $d(x, S)$ is the number of windows that are present in string x but not in any string in S . We are interested in exploring schemes that are intermediate between conjunctive and disjunctive detection (both positive or negative). In these intermediate schemes, a string may be treated as self, even if not all its window patterns have been observed before, i.e., $\text{Scheme}_M = \{x | d(x, S) \leq \tau\}$, where a τ is some threshold. Intuitively, strings that differ a small amount from the sample are more likely to be a part of RS than those that differ a lot. Extending this idea one step further, one can imagine distance measures that do not assign a uniform value to every match but instead take into account structural and statistical properties of the sample in order to weigh the relative merits of distinct matches.

E. Implementation Issues

Although we have emphasized the representational power of different detection schemes in this paper, there are important

implementation considerations that affect the choice of a match rule and detection scheme. Here, we briefly discuss the r -chunks match rule and contrast it with rcb, although, as shown in Section IV-C, the class of languages recognized by them are distinct. Nevertheless, the comparison highlights some important properties of r -chunks. We have thus far described our detection schemes as a set of detectors, first because we find it useful for visualizing and analyzing their properties and, second, because we consider that some applications might benefit from distributing them throughout several nodes; thus, we make no assumptions as to their location. Nevertheless, it is important to note that there is a large body of work [67] regarding string matching that suggests a more efficient representation of detectors along with the algorithms that operate on them but impose additional restrictions on the ease of distributing them. With this in mind, we consider three implementation issues: the cost of generating detectors, the cost of storing detectors, and the cost of locating detectors.

1) *Detector Generation*: Given a self set S , a straightforward method for generating the appropriate r -chunks detectors is to search the entire sample to determine which patterns are present in each window, requiring $O(t|S|)$ time and $O(t2^r)$ space. Note that the detector set can be sorted during its generation at no extra cost (using bin sort and assuming $|S| > 2^r$). There have been several algorithms proposed for generating detectors for the rcb match rule [3], [13], [18], [19] that report linear generation time and $O(t^2 2^r)$ extra space.

2) *Detector Storage*: Intuitively, r -chunks requires more space than rcb since every rcb detector (l bits each) can be decomposed into t r -chunks detectors (tr bits). Nevertheless, depending on the particular structure of S , there is likely to be a significant number of repeated patterns in each window amongst the set of rcb detectors. Because r -chunks represents each distinct pattern once, this could lead to smaller space requirements in some cases. In general, the space required for r -chunks is $O(tr2^r)$.

3) *Detector Location*: Having detectors specific for each window (as in r -chunks) allows them to be stored in sorted order for each window. Thus, checking to see if a string belongs to self or not requires $O(tr)$ time, whereas inserting or deleting a detector is just $O(r)$. There is also a potential that hashing schemes could reduce the search time.

F. Crossover Closure and Genetic Algorithms

There is a tantalizing connection between the crossover closure described in this paper and the crossover operator often used in genetic algorithms. The crossover closure contains only a subset of all possible one-point crossovers in the self set, and therefore, the generalization achieved by r -contiguous bits is not identical with the space of coordinate hyperplanes defined by the crossover operator. Likewise, the r -chunks matching rule is nothing more than a restricted form of the schema notation often used to explain genetic algorithm behavior. An additional area of future investigation is to make these connections more precise and to determine if the Schema theorem from genetic algorithms bears any relation to the processing performed by either of the contiguous bits match rules.

VIII. CONCLUSION

In this paper, we presented a formal framework for analyzing different positive and negative detection schemes in the context of approximate matching. Although the framework we presented is quite general, the primary application we have in mind is anomaly intrusion detection, for example, detecting anomalous TCP connections in a local area network or detecting anomalous patterns of system calls in executing processes. We gave examples of how different partial matching rules fit into our scheme, including Hamming distance, r -contiguous bits, and n -grams. We characterized the generalization induced by r -contiguous bits (and its relative r -chunks) by defining the crossover closure of a sample of l -length strings. Next, we showed that the crossover closure is related to the concept of a lossless join in relational database theory.

With this formal apparatus in place, we were then able to give some theoretical results on the relative power of the different detection schemes and matching rules. In particular, we showed that the r -chunks match rule subsumes r -contiguous bits matching, that (under rcb or r -chunks) the class of languages recognized by negative disjunctive detection is the same as that of positive conjunctive detection, and that the crossover closure of a sample S exactly characterizes the class of languages recognized by negative disjunctive detection and positive conjunctive detection (under r -chunks matching).

Next, we considered the number of detectors that are required to provide maximal discrimination for fixed r under Scheme_{PC} and Scheme_{ND}. We gave closed-form expressions for both schemes, allowing us for the first time to estimate how large a self set must be before negative detection is a computationally advantageous strategy. We then discussed the expected size of the crossover closure and explored how the use of permutations of the representation can reduce the size of the crossover closure. Finally, we showed that the class of languages recognized by r -contiguous bits, augmented with permutation masks, properly contains that recognized by Hamming distance, and we showed that even with the use of permutations, there are some languages that cannot be recognized using r -contiguous bits matching.

The theoretical results presented here are significant because they address a large and quickly expanding body of experimental work in intrusion detection that uses one of the detection schemes. It is our hope that these theoretical results will make it easier to understand the experimental results and to predict which approaches are most promising for which problems.

ACKNOWLEDGMENT

The authors would also like to thank J. Holland, A. Tyrrell, W. Langdon, and the anonymous reviewers, who's comments greatly improved the presentation of this work.

REFERENCES

- [1] S. Forrest, A. S. Perelson, L. Allen, and C. R. Cheru R. Kuri, "Self-non-self discrimination in a computer," in *Proc. IEEE Symp. Research Security Privacy*, Los Alamitos, CA, 1994.
- [2] S. Forrest, S. Hofmeyr, A. Somayaji, and T. Longstaff, "A sense of self for Unix processes," in *Proc. IEEE Symp. Computer Security Privacy*, 1996.
- [3] P. D'haeseleer, S. Forrest, and P. Helman, "An immunological approach to change detection: algorithms, analysis and implications," in *Proc. IEEE Symp. Computer Security Privacy*, 1996.
- [4] D. Dasgupta and S. Forrest, "Novelty detection in time series data using ideas from immunology," in *Proc. Int. Conf. Intelligent Systems*, 1996.
- [5] S. Hofmeyr, A. Somayaji, and S. Forrest, "Intrusion detection using sequences of system calls," *J. Comput. Security*, vol. 6, pp. 151–180, 1998.
- [6] D. J. Smith, S. Forrest, D. H. Ackley, and A. S. Perelson, "Variable efficacy of repeated annual influenza vaccination," *Proc. National Acad. Sciences (PNAS)*, vol. 96, pp. 14 001–14 006, 1999.
- [7] S. Hofmeyr and S. Forrest, "Immunity by design: an artificial immune system," in *Proc. Genetic Evolutionary Computation Conf.*, San Francisco, CA, 1999, pp. 1289–1296.
- [8] —, "Architecture for an artificial immune system," *Evol. Comput. J.*, vol. 8, no. 4, pp. 443–473, 2000.
- [9] A. Somayaji and S. Forrest, "Automated response using system-call delays," in *Proc. Usenix Security Symp.*, 2000.
- [10] J. Balthrop, S. Forrest, and M. Glickman, "Revisiting LISYS: parameters and normal behavior," in *Proc. Congr. Evolutionary Computation*, 2002.
- [11] J. Balthrop, F. Esponda, S. Forrest, and M. Glickman, "Coverage and generalization in an artificial immune system," in *Proc. Genetic Evolutionary Computation Conf.*, 2002.
- [12] D. L. Chao and S. Forrest, "An aesthetic immune system (tentative title)," in *Proc. Artificial Life VIII: The 8th Int. Conf. Simulation Synthesis Living Systems*, May 2002.
- [13] P. D'haeseleer, "An immunological approach to change detection: theoretical results," in *Proc. 9th IEEE Computer Security Foundations Workshop*, 1996.
- [14] S. Hofmeyr, *An Immunological Model of Distributed Detection and its Application to Computer Security*. Albuquerque, NM: Univ. New Mexico Press, 1999.
- [15] P. D. Williams, K. P. Anchor, J. L. Bebo, G. H. Gunsch, and G. D. Lamont, "CDIS: Toward a computer immune system for detecting network intrusions," in *Proc. Fourth Int. Symp., Recent Adv. Intrusion Detection*, W. Lee, L. Me, and A. Wespi, Eds., Berlin, Germany, 2001, pp. 117–133.
- [16] J. Kim and P. J. Bentley, "An evaluation of negative selection in an artificial immune system for network intrusion detection," in *Proc. Genetic Evolutionary Computation Conf.*, San Francisco, CA, 2001, pp. 1330–1337.
- [17] S. T. Wierzchon, "Generating optimal repertoire of antibody strings in an artificial immune system," in *Intelligent Information Systems*, M. A. Klopotek and M. Michalewicz, Eds. Heidelberg, Germany, New York: Physica-Verlag, 2000, pp. 119–133.
- [18] —, "Discriminative power of the receptors activated by k -contiguous bits rule," *J. Comput. Sci. Technol.*, vol. 1, no. 3, pp. 1–13, 2000.
- [19] —, "Deriving concise description of nonself patterns in an artificial immune system," in *New Learning Paradigm in Soft Computing*, S. T. Wierzchon, L. C. Jain, and J. Kacprzyk, Eds. Heidelberg, Germany, New York: Physica-Verlag, 2001, pp. 438–458.
- [20] W. Lee and S. Stolfo, "Data mining approaches for intrusion detection," in *Proc. 7th USENIX Security Symp.*, 1998.
- [21] "An agent based architecture for a computer virus immune system," in *Proc. GECCO Workshop Artificial Immune Syst.*, D. Dasgupta, Ed., 2000.
- [22] G. B. Lamont, R. E. Marmelstein, and D. A. Van Veldhuizen, "A distributed architecture for a self-adaptive computer virus immune system," in *New Ideas in Optimization*. London, U.K.: McGraw-Hill, 1999, Advanced Topics in Computer Science Series, pp. 167–183.
- [23] M. Burgess, H. Haugerud, and S. Straumnes, *Measuring System Normality I: Scales and Characteristics*, 2000.
- [24] C. Marceau, "Characterizing the behavior of a program using multiple-length N -grams," in *Proc. New Security Paradigms Workshop*, Cork, Ireland, 2000.
- [25] M. J. Kearns and U. V. Vazirani, *An Introduction to Computational Learning Theory*. Cambridge, MA: MIT Press, 1994.
- [26] J. C. Schlimmer, *Concept Acquisition Through Representational Adjustment*. Irvine, CA: Univ. California Press, 1987.
- [27] D. Angulin, "Learning regular sets from queries and counterexamples," *Inform. Comput.*, vol. 75, pp. 87–106, 1987.
- [28] T. Lane, *Machine Learning Techniques for the Computer Security Domain of Anomaly Detection*. West Lafayette, IN: Purdue Univ. Press, Aug 2000.

- [29] Lincoln Laboratories. (1999) DARPA Intrusion Detection Evaluation. [Online] <http://www.ll.mit.edu/IST/ideval/index.html>
- [30] D. W. Bradley and A. M. Tyrrell, "The architecture for a hardware immune system," in *The Third NASA/DoD Workshop on Evolvable Hardware*, D. Keymeulen, A. Stoica, J. Lohn, and R. S. Zebulum, Eds. Long Beach, CA: IEEE Computer Society Press, July 12–14, 2001, pp. 193–200.
- [31] —, "A hardware immune system for benchmark state machine error detection," in *Proc. Congr. Evolutionary Computation*, Honolulu, HI, May 2002.
- [32] —, "Immunotronics: Novel finite state machine architectures with built in Self test using Self-Nonself differentiation," *IEEE Trans. Evol. Comput.*, vol. 6, pp. 227–238, June 2002.
- [33] D. Dasgupta and F. Gonzalez, "An immunity-based technique to characterize intrusions in computer networks," *IEEE Trans. Evol. Comput.*, vol. 6, June 2002.
- [34] S. Sathyanath and F. Sahin, "Artificial immune systems approach to a real time color image classification problem," in *Proc. IEEE Int. Conf. Syst., Man, Cybern.*, 2001.
- [35] D. Q. Naiman, "Statistical anomaly detection via HTTPD data analysis," *Comput. Statist. Data Anal.*, to be published.
- [36] K. Tan and R. Maxion, "Why 6? Defining the operational limits of stide, and anomaly-based intrusion detector," in *Proc. IEEE Symp. Security Privacy*. Oakland, CA: IEEE Press, 2002.
- [37] J. W. Kappler, N. Roehm, and P. Marrack, "T-cell tolerance by clonal elimination in the thymus," *Cell*, vol. 49, pp. 273–280, Apr. 24, 1987.
- [38] J. K. Percus, O. Percus, and A. S. Perelson, "Predicting the size of the antibody combining region from consideration of efficient self/nonself discrimination," in *Proc. Nat. Acad. Sci.*, vol. 90, 1993, pp. 1691–1695.
- [39] —, "Probability of self-nonself discrimination," in *Theoretical and Experimental Insights into Immunology*, A. S. Perelson and G. Weisbuch, Eds. New York: Springer-Verlag, 1992.
- [40] P. F. Brown, V. J. Della Pietra, P. V. de Souza, J. C. Lai, and R. L. Mercer, "Class-based n-gram models of natural language," in *Proc. IBM Natural Language ITL*, Paris, France, 1990.
- [41] D. E. Denning, "An intrusion-detection model," *IEEE Trans. Softw. Eng.*, vol. SE–2, p. 222, Feb. 1987.
- [42] H. Teng, K. Chen, and S. Lu, "Adaptive real-time anomaly detection using inductively generated sequential patterns," in *Proc. IEEE Symp. Research Computer Security Privacy*, Los Alamitos, CA, 1990.
- [43] H. S. Javitz and A. Valdes, "The SRI IDDES statistical anomaly detector," in *Proc. IEEE Symp. Research in Security and Privacy*, 1991, pp. 316–326.
- [44] H. S. Javitz, A. Valdes, T. F. Lunt, A. Tamaru, M. Tyson, and J. Lowrance, *Next Generation Intrusion Detection Expert System (NIDES)*. Menlo Park, CA: SRI Int., 1993.
- [45] P. Helman, G. Liepins, and W. Richards, "Foundations of intrusion detection," in *Proc. Fifth Computer Security Foundations Workshop*, Franconia, NH, 1992, pp. 114–120.
- [46] M. Roesch. SNORT. [Online] <http://www.snort.org/>
- [47] M. Damashek, "Gauging similarity with n-grams: Language-independent categorization of text," *Science*, vol. 267, pp. 843–848, 1995.
- [48] A. Arning, R. Agrawal, and P. Raghavan, "A linear method for deviation detection in large databases," in *Proc. 2nd Knowledge Discovery Data Mining*, 1996, pp. 164–169.
- [49] V. Barnett and T. Lewis, *Outliers in Statistical Data*. New York: Wiley, 1994.
- [50] E. Eskin, "Anomaly detection over noisy data using learned probability distributions," in *Proc. 17th Int. Conf. Machine Learning*, San Francisco, CA, 2000, pp. 255–262.
- [51] D. H. Fisher, "Knowledge acquisition via incremental conceptual clustering," *Machine Learning*, vol. 2, no. 2, pp. 139–172, 1987.
- [52] D. Freedman, R. Psani, and R. Purves, *Statistics*. New York: W. W. Norton, 1978.
- [53] H. Garcia-Molina, J. Ullman, and J. Widom, *Database Systems: The Complete Book*. Englewood Cliffs, NJ: Prentice-Hall, 2001.
- [54] S. J. Hanson and M. Bauer, "Conceptual clustering, categorization, and polymorphy," *Machine Learning*, vol. 3, no. 4, pp. 343–372, 1989.
- [55] P. Helman and G. Liepins, "Statistical foundations of audit trail analysis for the detection of computer misuse," *IEEE Trans. Softw. Eng.*, vol. 19, pp. 886–901, Sept. 1993.
- [56] P. Helman and J. Bhangoo, "A statistically based system for prioritizing information exploration under uncertainty," *IEEE Trans. Syst., Man, Cybern.*, vol. 27, pp. 449–466, July 1997.
- [57] P. Helman and R. Gore, "Prioritizing information for the discovery of phenomena," *J. Intell. Inform. Syst.*, vol. 11, no. 2, pp. 99–138, Sept./Oct. 1998.
- [58] D. Hoaglin, F. Mosteller, and J. Tukey, *Understanding Robust and Exploratory Data Analysis*. New York: Wiley, 1983.
- [59] A. Jain and R. Dubes, *Algorithms for Clustering Data*. Englewood Cliffs, NJ: Prentice-Hall, 1988.
- [60] L. Kaufman and P. Rousseeuw, *Finding Groups in Data*. New York: Wiley, 1990.
- [61] E. Knorr and R. Ng, "A unified approach for mining: Properties and computation," in *Proc. 3rd Knowledge Discovery Data Mining*, 1997, pp. 219–222.
- [62] —, "Algorithms for mining distance based outliers in large databases," in *Proc. 24th VLDB*, 1998, pp. 392–403.
- [63] D. Maier, *The Theory of Relational Databases*. New York: Computer Science Press, 1983.
- [64] R. Ng and J. Han, "Efficient and effective clustering methods for spacial datamining," in *Proc. 20th VLDB*, 1994, pp. 144–155.
- [65] G. Salton and M. J. McGill, *Introduction to Modern Retrieval*. New York: McGraw-Hill, 1983.
- [66] J. T. Tou and R. C. Gonzalez, *Pattern Recognition Principles*. Reading, MA: Addison Wesley, 1974.
- [67] D. Gusfield, *Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology*. Cambridge, U.K.: Cambridge Univ. Press, 1997.
- [68] B. Scholköpfung and A. Smola, *Learning with Kernels*. Cambridge, MA: MIT Press, 2002.
- [69] C. Carla Marceau. Characterizing the behavior of a program using multiple-length n-grams. presented at *Proc. New Security Paradigm Workshop*. [Online] <ftp://ftp.oracorp.com/documents/MultiLengthStrings.pdf>

Fernando Esponda received the degree in computer engineering from the Instituto Tecnológico Autónomo de México (ITAM), Mexico City, in 1995 and will soon receive the masters degree from the Instituto de Investigaciones en Matemáticas Aplicadas y Sistemas (IIMAS-UNAM), Mexico City. He is currently pursuing the Ph.D. degree at the University of New Mexico, Albuquerque. His research interests are in biologically inspired adaptive systems and machine learning.

Stephanie Forrest received the B.A. degree from St. John's College, Santa Fe, NM, and the M.S. and Ph.D. degrees from the University of Michigan, Ann Arbor.

She is currently a Professor of computer science at the University of New Mexico (UNM), Albuquerque, and a member of the residential faculty at the Santa Fe Institute. Before joining UNM, she was with Teknowledge Inc., Palo Alto, CA, and was a Director's Fellow at the Center for Nonlinear Studies, Los Alamos National Laboratory, Los Alamos, NM. Her research interests are in adaptive systems, including genetic algorithms, computational immunology, biological modeling, and computer security.

Paul Helman was born in Brooklyn, NY, in 1954. He received the B.A. degree from Dickinson College, Carlisle, PA, in 1976, the M.S. degree from Stanford University, Stanford, CA, in 1977, and the Ph.D. degree from the University of Michigan, Ann Arbor, in 1982.

Currently, he is a Professor of computer science at the University of New Mexico, Albuquerque. His research interests include bioinformatics, machine learning, data mining, database theory, and combinatorial optimization. He also has authored two computer science text books: (with R. Veroff) *Walls and Mirrors: Intermediate Problem Solving and Data Structures* (Menlo Park, CA: Benjamin Cummings, 1986) and *The Science of Database Management* (Burr Ridge, IL: Richard D. Irwin, 1994).