

Probabilistic anomaly detection in distributed computer networks

Mark Burgess

Oslo University College, Cort Adelers gate 30, 0254 Oslo, Norway

Abstract

Distributed host-based anomaly detection has not yet proven practical due to the excessive computational overhead during training and detection. This paper considers an efficient algorithm for detecting resource anomalies in event streams with either Poisson or long-tailed arrival processes. A form of distributed, lazy evaluation is presented, which uses a model for human-computer interaction based two-dimensional time and a geometrically declining memory to yield orders of magnitude improvements in memory requirements. A three tiered probabilistic method of classifying anomalous behaviour is discussed. This leads to a computationally and memory economic means of finding probable faults amongst the symptoms of network and system behaviour.

Key words: Machine learning, anomaly detection, data-mining.

1 Introduction

Computer anomaly detection is about discerning regular and irregular patterns of behaviour, in the variables that characterize computer systems. The detection of anomalies in computer systems has often been pursued as the unambiguous goal of searching for potential breaches of security; it goes often hand in hand with Network Intrusion Detection, in which content-analyses of data are performed in real time with the aim of revealing suspicious activity[1–5]. However, this is only one application for anomaly detection; computers can also be approached as self-regulating systems that respond to changes in their environment in order to stabilize their own behaviour. In that case, anomaly detection becomes an integral part of the system's regulatory process. Previously, the cost of performing such an analysis on every host has been prohibitive, but this paper will suggest a way of overcoming this difficulty.

Anomaly detectors apply machine learning and analysis to see whether any long term trends can be found in data. One such approach was suggested in the early 1990s and has recently been revived [6,7]. Automated self-regulation in host management has also been discussed in refs [8–10], as well as adaptive behaviour [11] and network intrusion detection [12,4]. Other authors have likened such mechanisms to immune systems, striking the analogy between computers and other collective systems in sociology and biology [13,14,10].

The ultimate aim of anomaly detection systems is to have adaptive behaviour that reponds in ‘real-time’, so that problematical events can be countered as quickly as possible. However, normal behaviour can only be determined by learning about past events: trends take time to learn and analyze. This paradox can only be resolved by modelling future behaviour, based on a statistical idealization of the past and an observation of the present (like weather forecasting). Even then, a timely response requires a rapid processing of observations. The computational burden of real-time anomaly detection can be considerable. One would thus like to spread the burden as far as possible around the network to minimize the load at any particular place.

This paper is motivated by two goals: to develop an efficient method of anomaly detection that avoids bottlenecks, and implements ‘lazy evaluation’ to avoid unnecessary computational burden; and to develop a language for expressing one’s *policy* about what constitutes an anomalous occurrence, relative to what has already been learned about the signal in the past. We shall make some progress towards both of these goals. The paper is organized as follows:

- (1) We begin with a brief summary of the idea of host based anomaly detection, its aims and motivations in relation to the future challenges of mobile and pervasive computing.
- (2) Existing techniques for mapping out empirical data characteristics are summarized and appropriate statistical measures for discussing normality are identified.
- (3) The notion of policy is then introduced, to account for the arbitrary aspects of data analysis, such as threshold values and the representation of corroborating environmental information that is not represented in the learning abilities of the nodes.
- (4) Based on the known characteristics of host data, a pseudo-periodic parameterization of time series is developed, which partitions the arrival process into weekly units. Some comments are made about data distributions and the implications for machine learning.
- (5) A description of the limited span, unsupervised learning algorithm, with predictable ‘forgetting power’, is presented.
- (6) Finally, a multi-stage classification of data is proposed, where a response is instigated only if a probabilistic detector signals a *probably significant* event (lazy evaluation).

2 Host-based observation of anomalies

In contemporary network design, traffic congestion is avoided by packet switching. i.e. by isolating traffic to ‘parallel’ branches of a network spanning tree that is absolutely necessary to deliver data to their destinations. Computers, or *nodes*, occupy points at the leaves of these branches and therefore experience an individual (subjective) view of the traffic.

Each computer in a network has a different experience of the environmental bath of requests and replies that commit its resources. Because the concept of an anomaly is a subjective one (what is unusual for one node is a regular occurrence for another), one might imagine that nowhere in the network is better equipped to reveal anomalies than the nodes at which they finally arrive. That is the viewpoint we shall adopt here, justified by a previous study of the same data discussed in this paper[15].

Traditionally, anomaly detection has been centralized to trunk limbs of a network, in the belief that one can only see the big picture if one is in possession of as much of the data as possible at a single place. Some approaches even attempt to combine the streams at different observation points into a single stream[16]. There are two problems with this strategy: it places the entire burden of analysis at a single location, and it does not gain access to anomalies that occur on non-network variables (e.g. disk usage), since these are never transmitted across the network.

The utility of centralized analysis cannot be completely dismissed, but it can be shown to have a limited value[15], and its clear disadvantage is the bottlenecking of traffic that is contrary to modern network design; placing the burden of computation at a single location, meaning that additional computing facilities are required. Studies, at Oslo University College, find that correlations between nodes are generally too noisy to be useful, unless there is an obvious functional relationship between hosts by design. The conclusion is that there is little to be gained by sharing resource data between hosts[15], since the only clear results are already known in the logistic map of services for the network. Hence there is no pressing need for centralization or serialization of traffic.

Another compelling reason for abandoning the idea of serialization of the data stream is that computers will soon be ubiquitous and devices will be transmitting and receiving data without any regard for a centralized authority, over unguided media. In such a world, the strategy of trying to centralize anomaly detection, at a single gateway, is flawed. A detection scheme, in which each host node is responsible for itself and no others, reflects the true distributed governance of the network and embodies the move from monolithic centralized

control to the more ‘free market economy’ approach.

The present work is based on the idea of computer immunology[17,10], in which one considers every computer to be an independent organism in a network ecology. The cfengine project places the individual computer rather than the network centre stage, in the belief that soon a majority of nodes will not be aligned with any centralized authority[8,10]. Serialization of traffic is abandoned and one combines the analysis of network variables with the analysis of internal host variables, using the natural filtration of data by packet switching.

Arriving events have to be classified and counted in order to identify their statistical significance. They have internal attributes, such as names, addresses, values etc. with semantic interpretations. These internal attributes contain information that can be used to identify what is meant by an anomaly, by placing events in a specific context and category. An anomaly engine is therefore a ‘prism’ and ‘spectrometer’, or a decision tree that expands an event arrival/renewal process[18] into a spectrum of distinct attributes. By looking at these attributes, with policy criteria that are appropriate for each, and then reassembling the information into a consistent picture, we perform something analogous to a ‘medical scan’ of the incoming event, which allows us to determine its significance to the system.

All scientific observations are facilitated in the context of a model, and one must therefore formulate such a model for classifying the observations as they arrive, in order to derive their meaning. The separation of scales is a crucial aspect of any model[19]; the distinction between trend and event plays a special role for anomaly detection[20]. Measurements of time-series autocorrelations show that significant variations are only observed in trends over times of the order of greater than 20 minutes in human-driven activity[21]. Since an anomaly response time can be up to 30 minutes in most systems, whether they depend on humans or automation, there is no point in labelling data much more extensively than this, even though many hundreds, even thousands, of individual events can occur per minute. The trends, which we must learn, in the data do not change as quickly as the microscopic details of the data stream and do not need to be sampled more often than a typical rate of change.

The philosophy in this work thus diverges from the strategy of examining every event exhaustively. A change-event only acquires importance if can be successfully decoupled from a known trend. We thus use a compromise between autocorrelation of numerical event scales and macroscopic level correlations, and split time into granules of five minutes. The data collector measures signals for a whole granule before deciding how it should respond to each coarse-grained event.

One ends up with a decision based on the following spectrum of attributes:

- The significance of the arrival time (the granule label).
- The significance of the arrival rate (number per granule), relative to a trend (average number per granule).
- The probable uncertainty of in the assumed trend (a specified number of standard deviations above or below the average trend).
- Entropy content of the distribution of symbol content within granules (described below).
- The symbolic content of specific attributes themselves, collected over a granule.

The memory required to implement this characterization is quite small: it comprises the space required to store each measurement granule, plus the space required to remember the significant attributes within the measured granule, for a finite number of granule labels.

The remainder of the paper considers how to rationally optimally compare incoming granule observations to a memory of what is statistically normal, using an economical method.

3 Entropy: symbolic uncertainty

Interpreting the meaning in a data stream is central to the problem of defining when an event is an anomaly. ‘Anomalouslyness’ is a subjective judgement, made within the context of past experience, and can be codified into a ‘policy’ about what is sufficiently anomalous to warrant a response. We can use numerical estimates of significance (statistics), but a complete analysis must also take into account the symbolic attributes of the arrivals too.

In order to define a practical policy for what anomalies are, one must have a straightforward classification of criteria. This is a challenge for the symbolic content, which can comprise many different data types, but we can use a simple information theoretical measure as a first approximation.

In random event processes we learn distributions of values for observations, which have characteristic shapes, and hence characteristic uncertainties. A sharp distribution about a given value means low uncertainty; a broad or flat distribution signifies a highly uncertain value. In either case, one must take the uncertainty in data into account before drawing conclusions about it. This means that one requires a measure of that uncertainty, in order to acknowledge it and use it.

The simplest gauges of a distribution are its *moments*, of which the variance (second moment) is the most well-known. The square-root of this (or standard

deviation σ) is a simple scale of the uncertainty. However, except in the case of Gaussian distributed data, there is no clear relationship between the uncertainty and σ . In digital distributions, i.e. histograms with few classes, the standard deviation is inaccurate and clumsy.

A very convenient measure of uncertainty, for any type of data, is the Shannon entropy[22,19,9]. In the theory of information, the Shannon entropy is a numerical characterization of a value distribution $p_i = n_i / \sum_i n_i$, based on frequency counts n_i of types i .

$$S = - \sum_{i=0}^C p_i \log p_i, \tag{1}$$

$0 \leq i \leq C$ run over the distinguishable classes of observation, n_i is the number of events of type i and p_i is the normalized ‘probability’ of measuring an event of type i in a similar signal.

The value of the entropy embodies the learning that has occurred about the random process and provides a convenient scalar measure for making policy about a random process. The entropy has a minimum value of zero, when all the observations are in a single class, and a maximum value of $S = \log \sum_i n_i$ when each class is occupied equally. This provides an adaptive, relative scale that can be applied to any interval of observation.

Example 1 *Consider the entropy of a particular attribute: the origin IP address of a data stream: in a network data stream, packets come from different IP addresses. In the following example, we see that traffic has arrived from five different IP addresses (intrepreted as different symbol classes), but predominantly from the first address. It forms a sharply focused distribution, as cfengine identifies:*

```

Frequency: 157.158.24.40    |*****+ (47/53)
Frequency: 80.203.17.11    |* (1/53)
Frequency: 66.196.72.28    |* (1/53)
Frequency: 80.202.77.107   |** (2/53)
Frequency: 80.213.238.106  |** (2/53)
-
Scaled entropy of addresses = 12.7 %
(Entropy = 0 for single source, 100 for flatly distributed source)

```

This is a low entropy distribution, because most of the symbols (in this case IP addresses) are of a single type. The alphabet for this comparison is learned in situ, and is the set of five Internet Protocol addresses listed above.

If we measure entropy as a percentage of the maximum attainable value, then it can itself be classified into high, middle and low, using arbitrary thresholds.

This may then be used as a filter in policy description. If we need a more specific characterization, there is always the standard deviation (the square root of the second moment of the distribution).

4 Lazy attribute extraction

‘False positives’ or ‘ghost anomalies’ are events where current algorithms find problems that are bogus; they are the familiar lament of anomaly detection designers. The dilemma is to know when an anomaly is ‘false’ or when an anomaly is uninteresting. However, false and uninteresting are two rather different criteria. To call an anomaly false is to assume that we have pre-decided a policy for what is truly an anomalous event and what is not. To call an anomaly interesting is to suggest either that a feature of the data is not only abnormal but highly unusual or that it is usual but not according to a recognizable pattern. Unfortunately, both of these criteria are matters of opinion rather than absolute measuring sticks. What is missing from most network anomaly detectors is an ability to express policy decisions about what is desirable and undesirable information.

In the present work, it is assumed that false anomalies occur for two main reasons:

- Because one attempts to classify data inappropriately (without a model).
- Because the policy for distinguishing anomalies is over-constrained.

The latter is often a byproduct of the security applications of anomaly detection: one is easily duped into overt ‘cold war’ paranoia that leads to an arms race of sensitivity: the desire to scrutinize every event.

Looking, as many have, to biological detection in the immune system [13,14,10] for inspiration, one finds an excellent yet imperfect system that is cheap to operate. Cost is important: an immune system that was so expensive that it has to kill us to keep us alive would be of little use. A host must continue with its primary function while detecting and responding to anomalies.

The biological immune system is a multi-tiered reactor with many levels of detection and a short memory. Organisms tolerate small amounts of harmful material, but mobilize countermeasures once they begin to do damage[23]. The key method by which the immune system avoids responding to false positives is the use of *costimulation*. A confirmation signal is required (like a dual key system) to set off an immune response. The system is lazy, in that it need not look for the confirmation signal unless the probability of an anomaly is already high.

For a computer detection scheme, we can use the same approach. First we look for a probably anomaly by comparing observation to learned experience. If the event looks probable, we can consider the evidence determined from a supporting semantic model. This reduces the amount of processing involved in detecting anomalous behaviour to an absolute minimum, by using ‘lazy evaluation’:

- (1) The system learns the normal state of activity on a host.
- (2) New events are considered anomalous if reliable data can place them at some sufficient number of standard deviations above the expected value at a given time of week.
- (3) The meaning of ‘sufficient’ must be defined as a matter of policy, in a given context. It is subjective.
- (4) If an event is found anomalous, it is dissected through our ‘prism’ in terms of its informational entropy and symbolic content.

The policy referred to here can be used to describe which anomalies are *interesting*, and specify when to respond.

Definition 1 (Anomaly identification policy) *An anomaly identification policy is a specification of predicates and thresholds, for observable attributes in a stream of observations, that is used to identify anomalies. It maps the set of observable attributes into the Boolean set $\{True, False\}$.*

The strategy used here is to base anomaly policy on what has been learned about past behaviour. For this, we measure the significance of events relative to known trends and variations, and then use the symbolic content of the events, if statistical anomalies are found, to determine how we should respond to them. This breakdown is important, because it emphasizes the need for a *policy* for describing the importance of events in a local environment. A policy codifies information that is not available by direct observation of the host state (information that would require evolutionary timescales to incorporate in biological systems) and is therefore an important supplement to the regulatory system.

Example 2 *For example, suppose one observes anomalous amounts of World Wide Web traffic often come from a single IP address source. Given no further information, one might dismiss a lot of traffic from a single source as a scan by an Internet search engine. However, search engines generally scan from a number of IP addresses in parallel. The IP address entropy of a friendly search engine scan is relatively high. So one might not be worried about high entropy, high traffic combinations.*

By examining the IP addresses contributing to a granule, and trying to resolve them, however, one sees that low entropy sources are sometimes associated with unregistered IP addresses (those which are not in the Domain Name Service

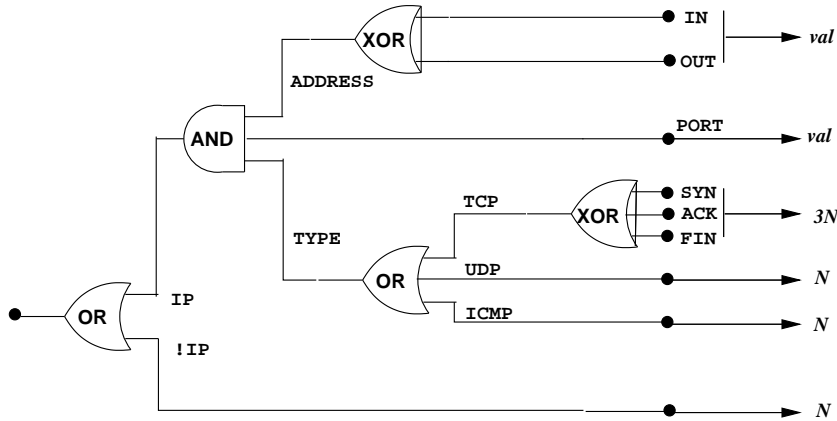


Fig. 1. An example network prism that splits an incoming event into generic categories. The signal enters at the left hand side and is classified as it passes to the right. This can be viewed as a reversed sequence of logic gates. One ends up with six frequency variables of magnitude N that count arrivals and two symbolic values.

or DNS). Such addresses are often hijacked or spoofed addresses and make one immediately suspicious of the source. Hence one can use this information and now codify a policy of responding to low entropy statistical anomalies from unregistered IP addresses. Policy is therefore a specification of acceptable attributes, here: arrival rate, address entropy and address resolvability.

Figure 1 shows an example of how one can easily split the example of a multifaceted network event into separate attributes that can be evaluated. The incoming packet is first examined to see if it is an IP (Internet Protocol) packet. If so, it has an address, a port number (except for ICMP) and a ‘layer 3’ encapsulation type (TCP,UDP etc). The different kinds of events can be counted to learn their statistical significance (we call these counting variables) and the remaining symbolic information (Internet addresses and port numbers) can be stored temporarily while the current sample is being analyzed. A sample is a coarse grained ensemble of events, collected over a five minute interval.

We now have two questions: how are data learned and how are events identified as statistically significant?

5 Arrival-renewal processes and self-similarity

A question that has been raised in recent years is that of the type of arrival or renewal process experienced by the end nodes in the network. This is often relevant for network analyses in which one attempts to model anomalies by looking at inter-arrival times of events, i.e. especially where one attempts to invoke memory of the recent past to track persistent events like connections.

Traditionally arrival processes have been assumed to be memoryless Poisson processes and analyses have used time correlations[24,2,25] to gauge likelihood of anomaly, but measurements of network traffic and indeed computer behaviour in general show that the arrival processes of normal computer operations sometimes have long tails and exhibit power law behaviour[26].

If a renewal process is *stable*, one has the chance of characterizing it with a time-series model. Two kinds of stable distributions are known for the inter-arrival times of renewal processes: Poisson distributions, and generalized stable Lévy process[27] (power laws). Aside from these, all other distributions are unstable under convolution[19].

Unfortunately, traffic volumes are rarely high enough to see stable distributions at leaf nodes. This has consequences for the analysis the time-series: statistical quantifiers sometimes diverge or become ill-defined. In particular, correlations over absolute time have little value.

The two inter-arrival time distributions of interest are the Poisson and the power-law:

$$N(\Delta t, n = 1) = (\lambda\Delta t)^1 e^{-\lambda\Delta t} = \sum_{m=0}^{\infty} \frac{(-1)^{-1}(\lambda\Delta t)^{m+1}}{m!} \quad (2)$$

$$N(\Delta t) = \mu(\Delta t)^{-\alpha}. \quad (3)$$

In general, both λ and μ can be functions of time, in the presence of trends. We shall make us of this below.

We can roughly gauge the type of process by an approximate measure of its degree of self-similarity, called the Hurst exponent H . This is a scaling exponent for the time-series over an range of average granule sizes. In other words, one assumes a general scaling law for a time-series observation $q(t)$ made at time t :

$$s^{-H} q(st) = q(t). \quad (4)$$

One then applies this to locally averaged functions:

$$s^{-H} \langle q(st) \rangle = \langle q(t) \rangle, \quad (5)$$

where $\langle \cdot \rangle$ is defined in eqn. 11. The exponent H can be estimated for real data by noting that, over an interval Δt ,

$$\langle \max(q(t)) - \min(q(t)) \rangle_{s\Delta t} = s^H \langle \max(q(t)) - \min(q(t)) \rangle_{\Delta t}, \quad (6)$$

i.e.

$$H = \frac{\log\left(\frac{\langle \max - \min \rangle_{s\Delta t}}{\langle \max - \min \rangle_{\Delta t}}\right)}{\log(s)}. \quad (7)$$

The values above unity signify probable power laws, with bursty behaviour. The data used in this paper fall into two main groupings. Some data for these are summarized by the in table 1. The results show a wide variety of

$q(\tau)$	$H(q)$
Users and processes	0.6 ± 0.07
Network connections (various)	$1.0 - 2.1 \pm 0.1$

Table 1

Approximate Hurst exponent ranges for different variables, show that the data exhibit variety of scaling behaviours once projected into a periodic framework. Some show long-tail indications while others have only Gaussian jitter.

behaviours in the signal, as measured over many months, some of which would tend to indicate self-similar behaviour. One therefore expects to have problems with the naive analysis of time correlations in these data. What is clear from comparing the actual graphs of data with their Hurst exponents is that the type of arrival process (characterized by H), is in no way correlated with the ability to separate signal from noise in this periodic parameterization discussed in this paper (characterized by small error bars, e.g. in fig. 3).

Another approach is required. In fact, we can avoid troubles associate with inter-arrival time tails entirely, below, by a simple transformation that integrates long or short inter-arrival time tails completely, by projecting them into a periodic time topology. This will leave us with a completely normalizable framework, with no ambiguities.

It was shown in ref. [28] that the transformation, to be described, may be used to represent the data in arrival processes, even with noise. Indeed, it is implicit in the Fourier theorem that any function projected into a periodic topology can be represented using a spectral (count per interval) representation.

6 Two-dimensional cyclic time parameterization

The solution to the arrival time issue is to make time itself finite and cyclic, using what is known about the processes in a computer system[29].

By making our model of time periodic we solve two problems: we are able replace the idea of time-correlation with a trend analysis based on counting,

i.e. a traditional frequency analysis, and we achieve a significant simplification of the machine-learning and analysis algorithms.

The basic observation that makes resource anomaly detection simpler and more efficient than traditional time-series analysis[30] is that there is an inhomogeneous pattern to human resource usage, and this is reflected in computer resource usage. This can be used to re-model and simplify the data. It yields, in effect, a simple and automatic supervised classification of the machine-learning process.

The approximate weekly periodicity observed in computer resources[29], allows one to parameterize time in topological slices of period $P = 1$ week, using the relation

$$t = nP + \tau, \quad (8)$$

$n = 0, 1, 2, \dots$ In this parameterization, time assumes a cylindrical form, labelled by two interleaved coordinates (τ, n) , both of which are discrete in practice[31]. We then sum or average over the n variable, leaving a single angular variable τ , after a renormalization of the distribution.

This transformation is indisputably justified in any process that is periodic in time, by Fourier theory, however, it may also be applied to functions that are only pseudo-periodic or close to periodic on average. See ref. [28] for a detailed discussion of this model. The key observation is that the arrival processes have periodic inhomogeneities $\lambda(\tau)$ and $\mu(\tau)$.

Applying this transformation to the problematical inter-arrival time distributions in eqn. (3), we see that the Poisson law may be written in terms of a periodic, dimensionless quantity τ/P , by substituting eqn. (8):

$$\begin{aligned} N(\Delta\tau/P, n) &= \sum_{m=0}^{\infty} \sum_{n=1}^{\infty} (-1)^m \frac{\lambda(\tau)^{m+n}}{n!m!} (nP + \Delta\tau)^{n+m} \\ &\propto f(\tau) \sum_{\alpha=-\infty}^{\infty} \zeta_{\alpha}(\Delta\tau/P) \end{aligned} \quad (9)$$

and the power law is simply:

$$N(\Delta\tau/P) = \mu(\tau) \sum_{n=0}^{\infty} (nP + \Delta\tau)^{-\alpha} \propto f'(\tau) \zeta_{\alpha}(\Delta\tau/P), \quad (10)$$

where f, f' are unknown, but periodic functions. Both of these results are expressed in terms of generalized forms known as zeta functions, multiplied by periodic amplitude functions. In this form, the Poisson law is seen to be

simply the generalization of the power law, formed by super-imposing many signals with different characteristic decay rates.

What is significant is that the resulting function of τ/P is well behaved for large times, over the limited interval $0 \leq \tau \leq P$; indeed, the large- n behaviour plays a less and less significant role that can be normalized away in the sums. The short τ behaviour is still singular for each power law component of positive α , but this is easily eliminated by coarse graining, which prevents $\delta\tau$ from ever being zero. These transformations are the standard tricks of renormalization statistics[32] (see section 8). To summarize, by restricting attention to fixed-size granules $\delta\tau$, taken one at a time, we can sum out and renormalize away the effects of the arrival process, leaving a periodic trend and an unknown amount of scatter.

As one would expect from the Fourier theory, a superposition of many such signals will lead to a strong trend if there are periodic variations in the data. Any other signals will appear as fluctuations that will either appear as noise or as anomalies, depending on their relative magnitude. In fig. 4, on the other hand, there is no convincing periodicity to be seen in the pattern of the data series, which begs the questions whether this method of insisting on periodicity is ‘appropriate’. However, the contention here is that this is not the way to look at it. The empirical studies indicate that weekly periodicity is, by far, the most important structural trend in human computer data[29]. If we can model the periodic parts of a signal, that are attributable to weekly behaviour, then it will be possible to extract those parts and get them under control. Everything that is left, relative to this trend, is either noise that can be renormalized away, or anomaly that can be identified more easily. Empirical studies provide compelling evidence for this.

We shall therefore *assume* that it is appropriate to project into periodic time, since the lack of periodicity is simply caused by a lack of a human-interaction and hence a lack of signal, meaning no pattern of normalcy and nothing would be gained by allowing time to extend indefinitely. In pedestrian terms, if we cannot separate significant variations from noise, in this parameterization, we cannot distinguish anomalies either, in this model, thus no harm is done by making the assumption.

The periodic parameterization of time means one can average (over n) the values at each point τ , leading to a mean and standard deviation of observations at each corresponding time of the week. Both the mean and standard deviations are thus functions of τ , and the latter plays the role of a scale for fluctuations at τ , which can be used to grade their significance. The cylindrical parameterization also enables one to invoke a compression algorithm on the data, so that one never needs to record more data points than exist within a single period. It thus becomes a far less resource intensive proposition to mon-

itor system normalcy. This compression is not possible in a linear time-series approach, using splines etc. [6].

Test data are taken to be a number of universal and easily measurable characters (see fig. 2):

- Number of users.
- Numbers of processes.
- Unix ‘load averages’.
- Average utilization of the system (load average).
- Number of incoming/outgoing connections to a variety of well know services.
- Numerical characteristics incoming and outgoing network packets.

These variables have been examined earlier and their behaviour is explained in[9,29]. Other variables might be studied in the future. A further advantage of this model is that the learned average behaviour can be stored indefinitely by using a simple database format, covering only a single working week in granules of five minutes, as we now show in the next section.

7 Separation of scales

In a dynamical, stochastic system, there are two basic kinds of change: non-equilibrium change (slow, progressive, trend variation that occurs on a timescale that is long compared to measurement) and fluctuations (occurring on a timescale that is fast compared to the measuring process). If the system is approximately stable, i.e. close to a steady state, then the combination of these can be used to characterize the recent history of the system. Fluctuations can be measured as a time series and analyzed[6] in order to provide the necessary information, and averaged out into *granules* or sampling intervals. During a sampling interval, data are collected, the mean and variance of the sample are found and these values are stored for the labelled interval. The sampling interval is chosen arbitrarily based on the typical auto-correlation length of the data being observed[29].

Time-series data can consume a lot of space. However, a considerable compression of the data can be achieved, and several orders of magnitude of computation time can be spared by separating the weekly data trends from the arrival of random events and by updating belief-estimates of only those trends, iteratively, rather than using an off-line analysis based on a complete journal of the past (see section 11).

This can render a good approximation to an appropriate sliding window, time-series data sample[29]. One obvious approach, for such a method, is to use a

convergent geometric series in order to define an average which degrades the importance of data over time: in other words, a series which forgets old data at a predictable rate. After a certain interval, the oldest memories in the data contribute only an insignificant fraction to the actual values.

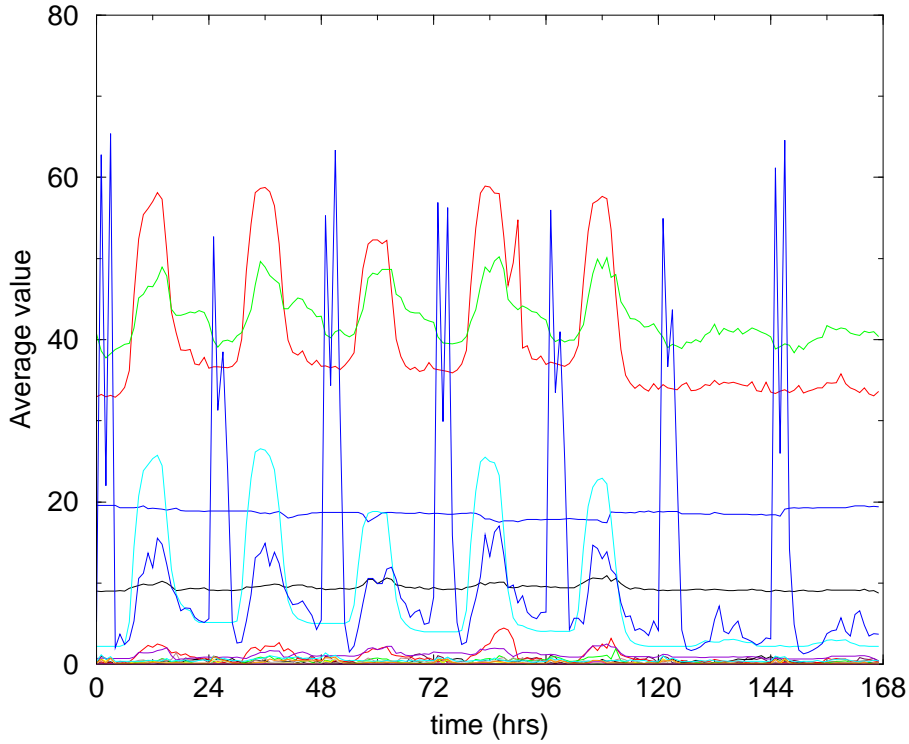


Fig. 2. A weekly periodogram of some resource variable expectations. These values are scaled and smoothed, measured by cfengine, and uncertainties have been suppressed. The lines represent the effective average thresholds for normal behaviour. Note how each line has its own characteristic ‘shape’ or pattern of usage, which the system learns by empirical measurement.

8 Computing expectations with memory loss

Our aim, then, is to realize the fact that we do not have to store the entire history of the system in order to infer its normal behaviour in the present. Rather, we can develop a Markov-style model in which the system not only learns but also forgets at a predictable rate.

Following the maintenance theorem of ref. [20], we define the normal behaviour of a system its *expected behaviour*. The standard deviation of the data *values* is a convenient scale by which to measure anomalies, and we now ignore the time-like nature of the arrival process for events.

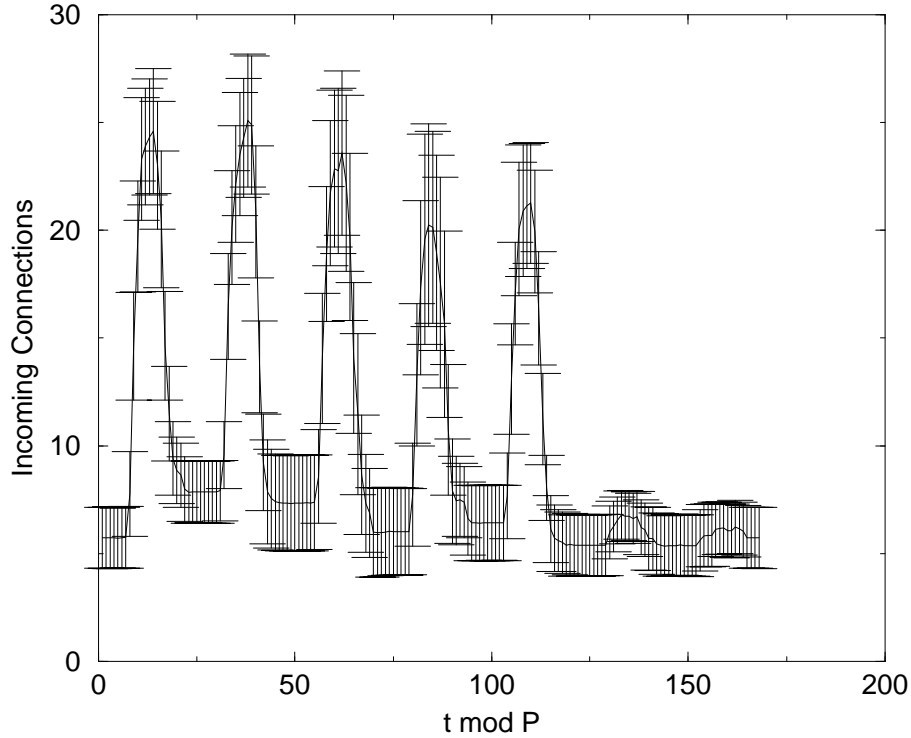


Fig. 3. Measured time trace of NETBIOS name lookups, averaged over 19.4 weeks. This basic pattern has been measured several times, starting with no data and has remained stable for almost two years. It has a clear signal. Uncertainties characterized by error bars represent the standard deviations $\sigma_{\langle P \rangle}(t \bmod P)$.

For a body of data, consisting of N data points $\{q_1, \dots, q_N\}$, one conventionally defines averages and standard deviations as follows.

$$\begin{aligned}
 \langle q \rangle &= \frac{1}{N} \sum_{i=1}^N q_i \\
 \langle q|Q \rangle &= \frac{1}{N} \sum_{i=1}^N q_i Q_i \\
 \sigma &= \sqrt{\frac{1}{N} \sum_{i=1}^n (q_i - \langle q \rangle)^2} \\
 &= \sqrt{\langle q^2 \rangle - \langle q \rangle^2} \\
 &= \sqrt{\langle \delta q | \delta q \rangle} \\
 &= \sqrt{\langle \delta q^2 \rangle}.
 \end{aligned} \tag{11}$$

This notation will help us to see that alternative definitions, with tailor-made properties, can unambiguously replace the expressions $\langle \cdot \rangle$ for averages. The use of these measures as characteristic scales in no way implies a model based

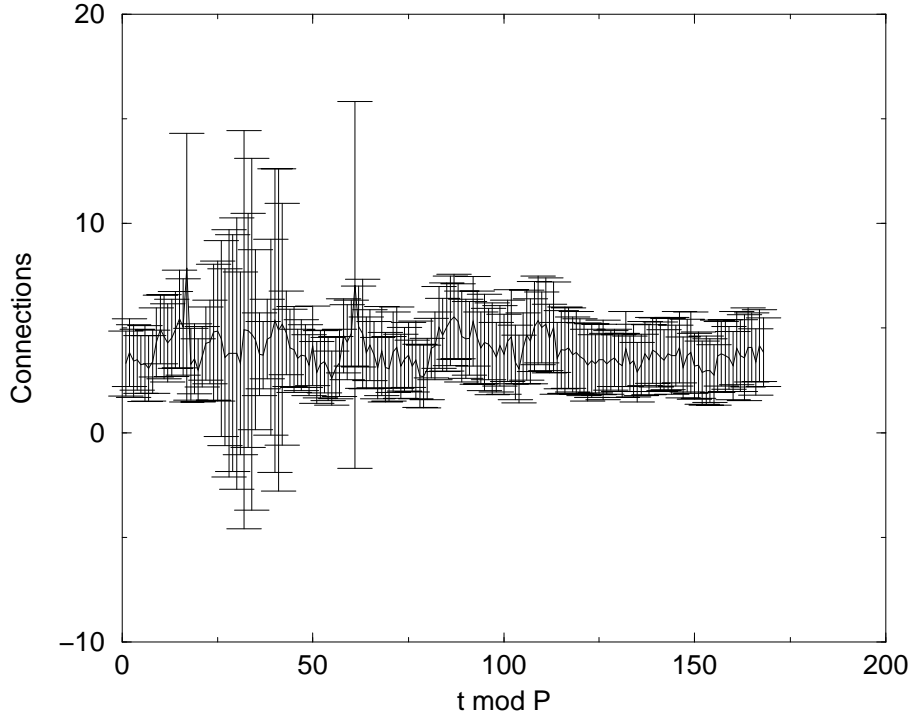


Fig. 4. Measured time trace averaged over 19.4 weeks of ftp connections. Here there is no discernable signal in the pattern of variations. The level of noise represented by the error bars is greater than the variation in the signal.

on Gaussian distributions.

To update the memory of averages and variances, an algorithm is required, satisfying the following properties:

- It should approximate an offline sliding-window time-series analysis that forgets old data at a predictable rate[29].
- It should present a minimal load to the system concerned.
- It must have a predictable error or uncertainty margin.

These goals can be accomplished straightforwardly as follows. We replace the usual expectation function with a new one with the desired properties, in such a way that derived quantities bear the same functional relationships as with the usual definitions.

$$\langle q \rangle \rightarrow \langle\langle q \rangle\rangle \quad (12)$$

that gradually forgets old data in a controlled manner. Similarly, we replace the standard deviation (or second moment of the data distribution) by

$$\sigma(\langle q \rangle) \rightarrow \sigma(\langle\langle q \rangle\rangle), \quad (13)$$

where

$$\sigma(\langle\langle q \rangle\rangle) \equiv \sqrt{[\langle\langle q^2 \rangle\rangle_N - \langle\langle q \rangle\rangle_N^2]} = \sqrt{\langle\langle \delta q^2 \rangle\rangle_N} \quad (14)$$

The new expectation function is defined iteratively, as follows:

Definition 2 (Iterative expectation function) *Let q be an observation, and $\langle\langle q_i \rangle\rangle$ be the i -th estimator of the average, with geometric fall-off, then $\langle\langle q_i \rangle\rangle$ may be defined by the recurrence relation:*

$$\begin{aligned} \langle\langle q \rangle\rangle_{i+1} &= (q | \langle\langle q \rangle\rangle_i) \\ \langle\langle q \rangle\rangle_0 &= 0. \end{aligned} \quad (15)$$

where

$$(q_1 | q_2) = \frac{w q_1 + \bar{w} q_2}{w + \bar{w}}. \quad (16)$$

and w, \bar{w} are constants.

Significantly, the number of data is now unspecified (we denote this by $i \rightarrow \infty$) meaning that this algorithm does not depend specifically on the arbitrary number of data samples N . Instead it depends on the ratio w/\bar{w} which is a forgetfulness parameter.

We note that, as new data points are measured after N samples, $\langle q \rangle$ changes only by q/N while $\langle\langle q \rangle\rangle_N$ changes by a fixed fraction $wq/(w + \bar{w})$ that is independent of N . Thus as the number of samples becomes large over time, the $\langle \cdot \rangle$ measure ceases to learn anything about the current state, as $q/N \rightarrow 0$, but $\langle\langle \cdot \rangle\rangle$ continues to refresh its knowledge of the recent past.

The repeated iteration of the expression for the finite-memory average leads to a geometric progression in the parameter $\lambda = \bar{w}/(w + \bar{w})$:

$$\begin{aligned} \langle\langle q \rangle\rangle_N \equiv (q_1 | (q_2 | \dots (q_r | (\dots | q_N)))) &= \frac{w}{w + \bar{w}} q_1 + \frac{\bar{w}w}{(w + \bar{w})^2} q_2 + \\ &\dots + \frac{w \bar{w}^{r-1}}{(w + \bar{w})^r} q_r + \dots + \frac{\bar{w}^n}{(w + \bar{w})^n} q_N. \end{aligned} \quad (17)$$

This has easily predictable properties. Thus on each iteration, the importance of previous contributions is degraded by λ . If we require a fixed window of size N iterations, then λ can be chosen in such a way that, after N iterations, the initial estimate q_N is so demoted as to be insignificant, at the level of accuracy

required. For instance, an order of magnitude drop within N steps means that $\lambda \sim |10^{-N}|$.

The learning procedure proposed here is somewhat reminiscent of a Bayesian probability flow[33,34], but it differs conceptually. A Bayesian algorithm assumes that each new datum can tell us the truth or falsity of a number of hypotheses. In our case, we have only single hypothesis: the normal state of the system, with a potentially unlimited amount of input. We do not expect this procedure to converge towards a static ‘true’ value as we might in a Bayesian hypothesis. Rather we want to implement a certain hysteresis in the normality function.

We now need to store the following triplets in a fixed-size database:

$$\{\tau, \langle\langle q \rangle\rangle(\tau), \sigma^2(\langle\langle q \rangle\rangle, \tau)\}. \quad (18)$$

We also use the δ symbol to represent the current deviation from average of a pseudo-periodic variable $q(t)$:

$$\delta q(t) \equiv q(t) - \langle\langle q \rangle\rangle_t \quad (19)$$

To satisfy the requirements of a decaying window average, with determined sensitivity $\alpha \sim 1/N$, we require,

- (1) $\frac{w}{w+\bar{w}} \sim \alpha$, or $w \sim \bar{w}/N$.
- (2) $\left(\frac{\bar{w}}{w+\bar{w}}\right)^N \ll \frac{1}{N}$, or $\bar{w}N \ll w$.

Consider the ansatz $w = 1 - r$, $\bar{w} = r$, and the accuracy α . We wish to solve

$$r^N = \alpha \quad (20)$$

for N . With $r = 0.6$, $\alpha = 0.01$, we have $N = 5.5$. Thus, if we consider the weekly update over 5 weeks (a month), then the importance of month old data will have fallen to one hundredth. This is a little too quick, since a month of fairly constant data is required to find a stable average. Taking $r = 0.7$, $\alpha = 0.01$, gives $N = 13$. Based on experience with offline analysis and field testing, this is a reasonable arbitrary value to choose.

9 Pseudo-periodic expectation

The recent behaviour of a computer can be summarized by n th order Markov processes, during periods of change, and by hidden Markov models during

steady state behaviour, but one still requires a parameterization for data points. Such models must be formulated on a periodic background[28], owing the importance of periodic behaviour of users. The precise algorithm for averaging and local coarse-graining is somewhat subtle, and involves naturally orthogonal time dimensions which are extracted from the coding of the database. It is discussed here using an ergodic principle: a bi-dimensional smoothing is implemented, allowing twice the support normally possible for the average, given a number of data points. This provides good security against “false positive” anomalies and other noise.

Consider a pseudo-periodic function, with pseudo-period P ,

$$\begin{aligned} q(t) &= \sum_{n=0}^{\infty} q(nP + \tau) & (0 \leq \tau < P) \\ &\equiv \sum_{n=0}^{\infty} \chi_n(\tau). \end{aligned} \quad (21)$$

This defines a set of periodic functions $\chi_n(\tau)$ with periodic coordinate $0 \leq \tau < P$. The time coordinate τ lives on the circular dimension. In practice, it is measured in p discrete time-intervals $\tau = \{\tau_1, \tau_2, \dots, \tau - p\}$. In this decomposition, time is a two-dimensional quantity. There are thus two kinds of average which can be computed: average over corresponding times in different periods (topological average $\langle \chi(\tau) \rangle_T$), and average of neighbouring times in a single period (local average $\langle \chi(\tau) \rangle_P$). For clarity, both traditional averages and iterative averages will be defined explicitly. Using traditional formulae, one defines the two types of mean value by:

$$\begin{aligned} \langle \chi \rangle_T(\tau) &\equiv \frac{1}{T} \sum_{n=l}^{l+T} \chi_n(\tau) \\ \langle \chi \rangle_P(n) &\equiv \frac{1}{P} \sum_{\ell=\tau}^{\tau+P} \chi_n(\ell) \end{aligned} \quad (22)$$

where l and τ are arbitrary start values and P, T are integer intervals for the averages, in the two time-like directions. Within each interval that defines an average, there is a corresponding definition of the variation and standard deviation, at a point τ :

$$\begin{aligned} \sigma_T(\tau) &\equiv \sqrt{\frac{1}{T} \sum_{n=l}^{n=l+T} (\chi_n(\tau) - \langle \chi \rangle_T(\tau))^2} = \sqrt{\langle \delta \chi_T | \delta \chi_T \rangle_T} \\ \sigma_P(n) &\equiv \sqrt{\frac{1}{P} \sum_{\ell=\tau}^{\ell=\tau+P} (\chi_n(\ell) - \langle \chi \rangle_P(\ell))^2} = \sqrt{\langle \delta \chi_P | \delta \chi_P \rangle_P}. \end{aligned} \quad (23)$$

Limited memory versions of these may also be defined, straightforwardly from the preceding section by replacing $\langle \delta q | \delta q \rangle$ with $\langle\langle \delta q^2 \rangle\rangle$ from eqn. (15)

$$\begin{aligned}\langle \chi \rangle_P &\rightarrow \langle\langle \chi \rangle\rangle_P \\ \langle \chi \rangle_T &\rightarrow \langle\langle \chi \rangle\rangle_T.\end{aligned}\tag{24}$$

Similarly, the deviations are given by

$$\begin{aligned}\sigma_{\langle\langle T \rangle\rangle}(\tau) &\equiv \sqrt{\langle\langle (\delta_{\langle\langle T \rangle\rangle} \chi)^2 \rangle\rangle_T} \\ \sigma_{\langle\langle P \rangle\rangle}(n) &\equiv \sqrt{\langle\langle (\delta_{\langle\langle P \rangle\rangle} \chi)^2 \rangle\rangle_P}\end{aligned}\tag{25}$$

where, for any measure X , we have defined:

$$(\delta_{\langle\langle P \rangle\rangle} X) \equiv X - \langle\langle X \rangle\rangle_P$$

$$(\delta_{\langle\langle T \rangle\rangle} X) \equiv X - \langle\langle X \rangle\rangle_T$$

(26)

(27)

Here one simply replaces the evenly weighted sum over the entire history, with an iteratively weighted sum that falls off with geometric degradation.

A major advantage of this formulation is that one only needs to retain and update two values per variable, the mean and the variance, in order to obtain all the information, not $2N$ data, for history size N .

10 Cross-check calibration: annealing potential anomalies

We now have a stable characterization of the time series that makes optimum use of the known structure of the data. In a two dimensional time series, one has two independent vectors for change that must be considered in generating a normal surface potential for comparison.

So far, the discussion has focused on a single periodicity in the time-series data, however we must also acknowledge the existence of sub-patterns within a single period. These patterns are not clear harmonics of the period, so they cannot be eliminated by redefinition of the period itself. Rather, they lead to apparent short term variations that, together with noise, can lead to apparent anomalies that are false.

It comes as no surprise to learn that the major sub-pattern is a daily one, once again driven by the daily 24 hour rhythm of activity, but it is not immediately clear why it is not the fundamental period of the system. The weekly

pattern can be reproduced with very low levels of noise, because the variations over many weeks of the weekly pattern are small. The daily pattern has much higher levels of uncertainty, since not all days are equivalent: weekends typically show very low activity and artificially increase the uncertainty in the expected signal. The difference between a weekend day and the variation in any day of the week over several weeks is significant, hence the working week yields the cleanest periodicity, at least in the data that have been collected in the present investigations[29].

One might perhaps expect that, in a clearly periodic signal, minor sub-patterns in observations would average out leaving a clear and smooth trend. This would render the problem of false anomalies insignificant, however, the favoured sensitivity of the new expectation function to recent events can also lead to artificial uncertainty. Random fluctuations at closely neighbouring times of day can also lead to apparent variations in the expectation function that is not statistically significant. We therefore define a procedure of smoothing, or annealing the anomalies by computing a *local average* as the smoothed vicinity of the current period. A traditional expectation expression for this would be:

$$\langle \chi \rangle_L(\tau) \equiv \frac{1}{L} \sum_{\ell=\tau-L/2}^{\tau+L/2} \langle \chi \rangle_T(\tau), \quad (28)$$

and in limited memory form, one has:

$$\langle \chi \rangle_L(\tau) \equiv \langle \langle \chi \rangle_T(\tau) \rangle_L, \quad (29)$$

and

$$\delta_{\langle L \rangle} \chi(\tau) \equiv \langle \chi \rangle_T - \langle \chi \rangle_L, \quad (30)$$

with corresponding measures for the standard deviations. Using these averages and deviation criteria, we have a two-dimensional criterion for normalcy, which serves as a control at two different time-scales. One thus defines normal behaviour as

$$\{\delta_{\langle L \rangle} \chi(\tau), \delta_P \chi(n)\} < \{2\sigma_{\langle L \rangle}(\tau), 2\sigma_{\langle P \rangle}(n)\}. \quad (31)$$

These may be simply expressed in geometrical, dimensionless form

$$\Delta(\tau, n) = \sqrt{\left(\frac{\delta_{\langle L \rangle} \chi(\tau)}{\sigma_{\langle L \rangle}(\tau)}\right)^2 + \left(\frac{\delta_{\langle P \rangle} \chi(n)}{\sigma_{\langle P \rangle}(n)}\right)^2}, \quad (32)$$

and we may classify the deviations accordingly into concentric, elliptical regions:

$$\Delta(\tau, n) < \begin{cases} \sqrt{2} \\ 2\sqrt{2} \\ 3\sqrt{2} \end{cases}, \quad (33)$$

for all τ, n . which indicate the severity of the deviation, in this parameterization. This is the form used by cfengine's environment engine[35]. Put simply, cfengine smudges gradients to make them smoother, as well as periodicities; hence, the smoothing algorithm is fully two-dimensional.

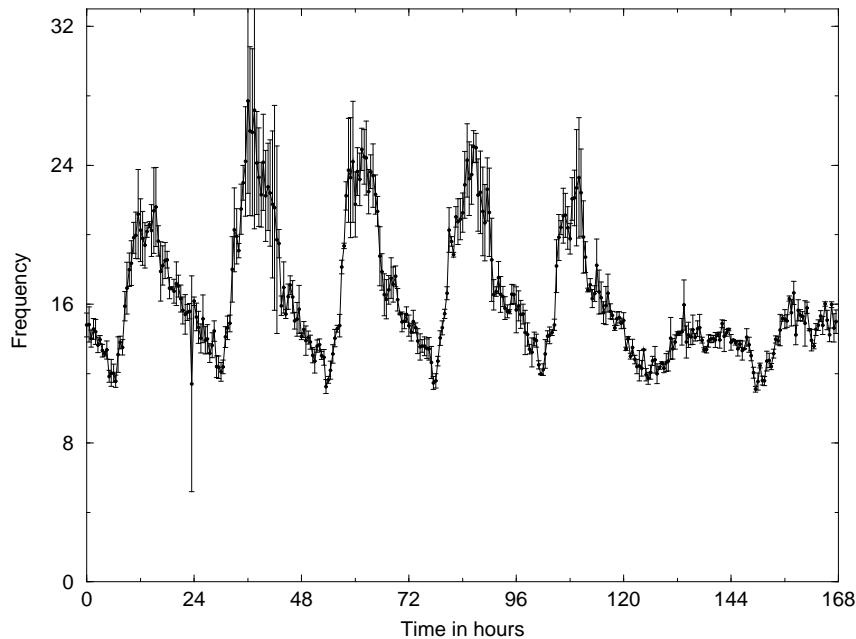


Fig. 5. A high resolution plot of process count behaviour, using the time-series sliding window from ref. [29]. This may be taken as a reference point. Data size is 21114875 bytes.

Figure 5 shows detailed, high resolution averages for process count data on a given host, over an interval of two months. At this level of detail one sees a jagged curve with error bars which vary considerably in magnitude. Figure 6 shows an image created using the iterative algorithm, on the same data, at the same resolution. The main trends of the curve are still visible. Error bars have been suppressed to avoid clutter (see fig. 7 for distribution characteristics).

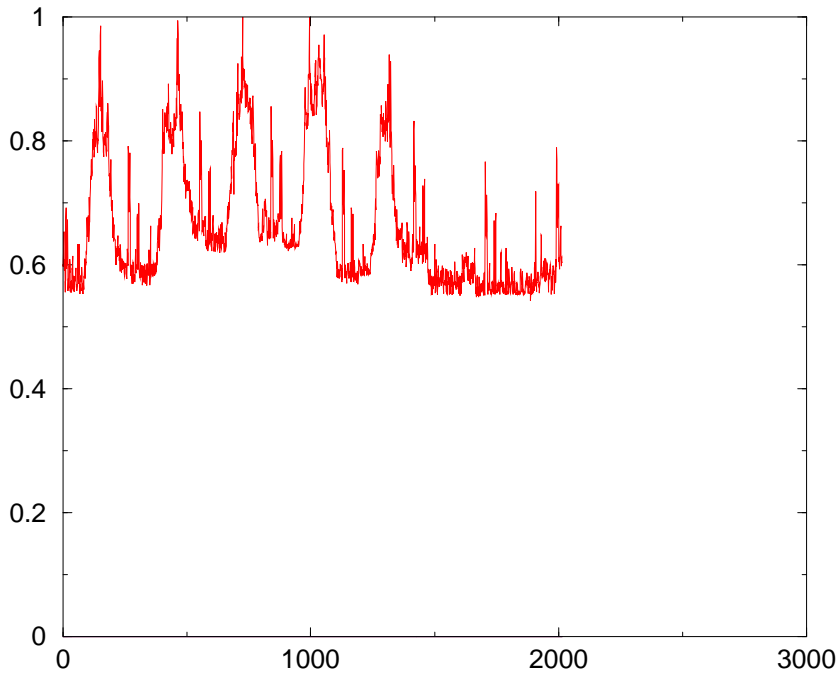


Fig. 6. A high resolution plot using the iterative algorithm using the the same process count data as in figure 5. Data size is 1531904 bytes (13.8 times smaller).

11 Co-stimulation: judging the calibrated sensitivity

The greatest problem for anomaly detection is in dealing with periods of low activity, i.e. a lack of recent past experience. In a period of low activity, every event is “abnormal”. Cold-spots pull the averages down and over-sensitize the detection of random events. In such a case, only a policy decision can renormalize the threshold level to avoid the detection of certain events.

Our anomaly measurement scale is the standard deviation, but a single standard deviation is often not even resolvable on a lightly used host, i.e. it is less than the discrete counting scale of the events; the appearance of a single new event might trigger a standard deviation from the norm. On more heavily loaded hosts, with persistent loading, more reliable measures of normality can be obtained.

How can one avoid a deluge of ‘false positives’ in anomaly detection? We must invoke policy to further classify events as interesting or uninteresting, using the information content of the events. As part of policy, we can combine the symbolic and numerical classifications of events in a kind of ‘co-stimulation’. In other words, one can use qualifiers (see fig. 8) to decide when to respond to the classified anomaly.

The scheme presented in this work has been implemented and tested as part of the cfengine project[35]. Cfengine is a distributed system administration agent,

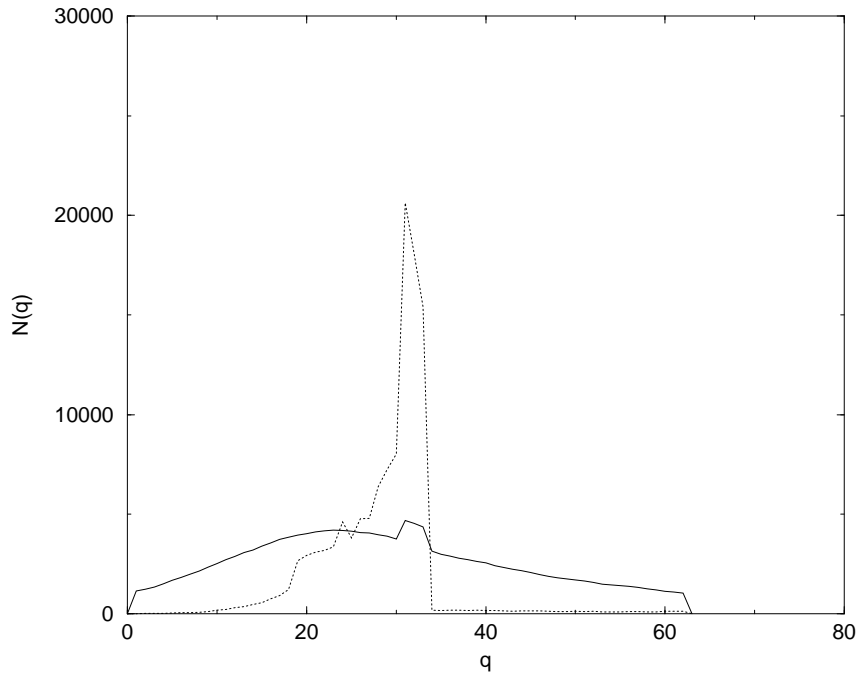


Fig. 7. The scaled scatter distributions, relative to the periodic trend, for HTTP traffic to and from a host. The solid line shows arriving connections, and the broken line shows outgoing connections. The peaks represent normal or expected value. The scatter of two values about normal is quite different, and far from Gaussian in the one case, and thus our idea of anomalous must also be different. The shapes of some of these curves can be calculated using the periodic model, see ref. [28] for details.

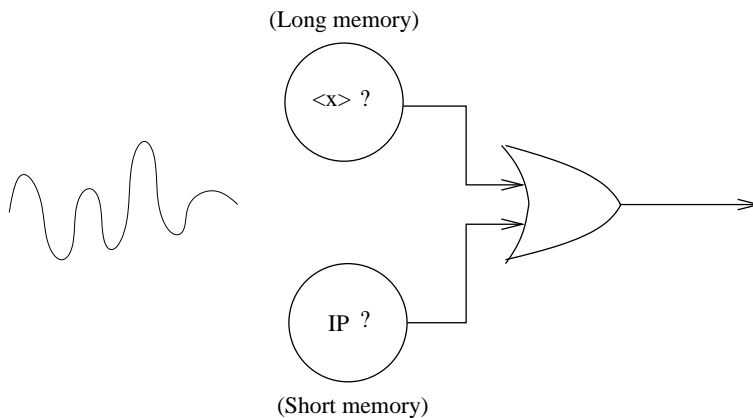


Fig. 8. A strategy of co-stimulation is used to sequentially filter information. First, long term (low grade) memory decides whether an event seems statistically significant and assess the likelihood of danger. If significant, short term (high grade) memory is used to recognize the source of the anomaly.

based loosely on the idea of a computer immune system. Anomaly detection is used to identify unusual behaviour that can diagnose problems of configuration and perhaps security. Cfengine already has a declarative language that can be based on the *classification of observations*. It serves as a useful test-bed for

anomaly policies.

Cfengine adds descriptive predicates for every true classification of its observations. For example, to generate a simple alert, one could write:

```
alerts:
```

```
entropy_smtp_in_low & anomaly_hosts.smtp_in_high_dev2::  
  
    "LOW ENTROPY smtp anomaly on $(host)/$(env_time)"  
  
    ShowState(incoming.smtp)
```

This would generate an alert if incoming E-mail (SMTP connections) exceeded two standard deviations above normal, at any given time of day, *and* if the traffic was predominantly from a single source (low entropy). Such an event could be a candidate for a ‘spam’ or junk-mail attack, for instance. The choice of how many standard deviations above expectation is a matter for policy, but it can be guided by the learned shapes of the data, such as those in fig. 7.

With this policy, a few anomalies per day is reasonable to expect from a moderately loaded host. At this level, anomaly warnings can maintain the attention of a human operator. Less significant anomalies can be handled in silence, by attaching programmed responses to the conditions that arise.

12 Some related work

The model used, in this work, to evaluate anomalies is unlike that of any previous systems known to the author. Nevertheless, it is fitting to place it in the context of other work. The literature on anomaly detection is vast and spans a variety of approaches and architectures, too numerous to chronicle. Only a few milestones are therefore mentioned below.

It is conventional, in the professional literature, to distinguish anomaly detectors from rule-based pattern matchers[12,36]. An anomaly detector uses the idea of unsupervised learning[37,38], whereas an intrusion detection system is based on a supervised expert database of already-identified patterns[39,40].

Most anomaly detectors attempt to learn temporal sequences, i.e. patterns or shapes that can be seen in a multivariate data stream. Early work in this for resource anomalies was carried out in [6], using time-series analysis and spline identification. Since then, more attention has been given to architectures for data collection, rather than for data analysis, even now e.g. [41]. This seems

to reflect a philosophy that more data will lead to a greater chance of finding anomalies, which is the apogee of the approach considered here. Today, packet analysis using the Snort[36] software is popular owing to its ready availability.

An important landmark in intrusion detection history was the SRI International (formerly Stanford Research Institute) Emerald project[24,42]. The aim of Emerald was to devise a layered approach to real-time anomaly detection in service traffic that was distributed and reintegrated centrally. Emerald was more of a collection system than an approach to anomaly decipherment. As such, the approach used is not necessarily at odds with that presented here, but it does not make apparent use of any compression based on observations of the data's structure. Its handling of statistics is somewhat unclear in the literature, but it was noted that it does not attempt to reduce false positives[43].

A more interesting technique has been used by Forrest et al., applying a probabilistic detection method, which is also inspired by immunology. The approach is fundamentally different to the work here: the authors employ stochastic searching of symbolic information streams, most notably digital streams of system codes. Notable references are [44–46], which inspired an application of this paper's methods as a supplementary enhancement in ref. [47]. No comparison of this method with the present one is possible in the present context. See ref. [47] for more details.

Statistical approaches have been employed in various refs. [48–50], but these are not like periodic trend compression seems to have been used.

The work that comes closest to the present approach is the Bayesian learning algorithm in ref. [51]. This author stops short of finding a way of using significant trends in the data to eliminate noise. In Bayesian learning, a system updates its belief about what to expect[37,52].

The signal analysis approach in ref. [53] takes an approach that is distantly related to the strategy, used in this paper, of using periodic functions to approximate the reconstituted time-series. It applies wavelet analysis to an unstructured time-series; this approach has become popular in the analysis of self-similar time-series. However, the analysis is an offline technique and did not yield clear advantages over cheaper methods.

The need for a policy to resolve subjective ambiguities in computer systems has been explored in a variety of access-security related contexts[54,55], but this is not a concept that has been discussed for pattern recognition. For intrusion and anomaly systems, policy usually amounts to defining lists of regular expressions to match symbolic traffic payloads. Although not all symbolic languages are regular, any finite symbolic language is regular[56], and all sequences are finite in practice. The computational simplicity of using regular expressions makes this approach the overwhelming approach of choice.

Policy is normally only applied to Intrusion Detection Systems and firewalls, rather than anomaly detection systems, see for example ref. [57]. Approaches that attempt to characterize and utilize the shape of statistical distributions, other than implicitly with a Gaussian model, are unknown to the present author.

13 Conclusions

The lazy host-based method of anomaly detection, used by cfengine, employs a two-dimensional time slice approach and a policy identification language. The resources required to store learned data are several orders of magnitude less than with traditional sliding window methods, due to an iterative learning scheme based on geometric series. It weights events on a sliding scale of importance so that recent events are geometrically more important than old events.

The periodic counting parameterization avoids problems associated with long tailed time-correlation divergences. It is unnecessary to assume Gaussian or Poisson statistics in their value and time distributions. The result is a probabilistic method, for detecting anomalous behaviour, that uses a policy-specified deviation from statistical expectation as a first sign of danger, and only then allows symbolic content to confirm or revoke the sign (co-stimulation).

The final aim of anomaly research is to have a turn-key, plug and play solution to the problem of detection, into which users insert their policy requirements and the machine does the rest. This requires a language for expressing anomalies in intuitive but *quantitative* terms. The model that is developed here allows policy to be expressed in terms of two statistical quantifiers:

- Number of arrival deviations above or below average, at any moment.
- The sharpness or bluntness of internal information within an observation.

Put in more information theoretical terms, this is nothing more than

- The statistical uncertainty of the environment.
- The symbolic uncertainty of the event.

With just these two parameters, the cfengine project has shown that one can elegantly classify and filter anomalies without supervision, and only a pre-specification of policy. Thus, we are able to say when an event is anomalous with surprising economy of expression and resources.

What it is still missing, in the technology of system regulation, is the determination of meaning in an anomaly, in the sense of how to mount a response

to it. This suggests finding an efficient characterization of the *non-statistical*, symbolic attributes of events, summarized in a *response policy* that attaches counter-measures to the anomalies. This is a separate and challenging problem which must be the subject of future work.

References

- [1] D. Denning. An intrusion detection model. *IEEE Transactions on Software Engineering*, 13:222, 1987.
- [2] V. Paxson. Bro: A system for detecting network intruders in real time. *Proceedings of the 7th security symposium. (USENIX Association: Berkeley, CA)*, 1998.
- [3] S. Forrest, S. A. Hofmeyr, A. Somayaji, and T. A. Longstaff. In *Proceedings of 1996 IEEE Symposium on Computer Security and Privacy (1996)*.
- [4] S. A. Hofmeyr, A. Somayaji, and S. Forrest. Intrusion detection using sequences of system calls. *Journal of Computer Security*, 6:151–180, 1998.
- [5] C. Kruegel and G. Vigna. Anomaly Detection of Web-based Attacks. In *Proceedings of the 10th ACM Conference on Computer and Communication Security (CCS '03)*, pages 251–261, Washington, DC, October 2003. ACM Press.
- [6] P. Hoogenboom and J. Lepreau. Computer system performance problem detection using time series models. *Proceedings of the USENIX Technical Conference, (USENIX Association: Berkeley, CA)*, page 15, 1993.
- [7] Y. Diao, J.L. Hellerstein, and S. Parekh. Optimizing quality of service using fuzzy control. *IFIP/IEEE 13th International Workshop on Distributed Systems: Operations and Management (DSOM 2002)*, page 42, 2002.
- [8] M. Burgess. A site configuration engine. *Computing systems (MIT Press: Cambridge MA)*, 8:309, 1995.
- [9] M. Burgess. Automated system administration with feedback regulation. *Software practice and experience*, 28:1519, 1998.
- [10] M. Burgess. Computer immunology. *Proceedings of the Twelfth Systems Administration Conference (LISA XII) (USENIX Association: Berkeley, CA)*, page 283, 1998.
- [11] M.I. Seltzer and C. Small. Self-monitoring and self-adapting operating systems. *Proceedings of the Sixth workshop on Hot Topics in Operating Systems, Cape Cod, Massachusetts, USA. IEEE Computer Society Press*, 1997.
- [12] M.J. Ranum et al. Implementing a generalized tool for network monitoring. *Proceedings of the Eleventh Systems Administration Conference (LISA XI) (USENIX Association: Berkeley, CA)*, page 1, 1997.

- [13] J.O. Kephart. A biologically inspired immune system for computers. *Proceedings of the Fourth International Workshop on the Synthesis and Simulation of Living Systems*. MIT Press, Cambridge MA., page 130, 1994.
- [14] S. Forrest, S. Hofmeyr, and A. Somayaji. *Communications of the ACM*, **40**:88, 1997.
- [15] K. Begnum and M. Burgess. Principle components and importance ranking of distributed anomalies. *Machine Learning Journal*, 58:217–230, 2005.
- [16] IETF. Intrusion detection exchange format.
- [17] A. Somayaji, S. Hofmeyr, and S. Forrest. Principles of a computer immune system. *New Security Paradigms Workshop, ACM*, September 1997:75–82.
- [18] G.R. Grimmett and D.R. Stirzaker. *Probability and random processes (3rd edition)*. Oxford scientific publications, Oxford, 2001.
- [19] M. Burgess. *Analytical Network and System Administration — Managing Human-Computer Systems*. J. Wiley & Sons, Chichester, 2004.
- [20] M. Burgess. On the theory of system administration. *Science of Computer Programming*, 49:1, 2003.
- [21] M. Burgess, G. Canright, and K. Engø. A graph theoretical model of computer security: from file access to social engineering. *International Journal of Information Security*, (accepted), 2004.
- [22] C.E. Shannon and W. Weaver. *The mathematical theory of communication*. University of Illinois Press, Urbana, 1949.
- [23] P. Matzinger. Tolerance, danger and the extended family. *Annu. Rev. Immun.*, 12:991, 1994.
- [24] H.S. Javitz and A. Valdes. The SRI IDES Statistical Anomaly Detector. In *Proceedings of the IEEE Symposium on Security and Privacy, May 1991*. IEEE Press, 1991.
- [25] V. Paxson and S. Floyd. Wide area traffic: the failure of poisson modelling. *IEEE/ACM Transactions on networking*, 3(3):226, 1995.
- [26] W.E. Leland, M. Taqqu, W. Willinger, and D. Wilson. On the self-similar nature of ethernet traffic. *IEEE/ACM Transactions on Networking*, pages 1–15, 1994.
- [27] K.I. Sato. *Levy Processes and Infinitely Divisible Distributions*. Cambridge studies in advanced mathematics, Cambridge, 1999.
- [28] M. Burgess. The kinematics of distributed computer transactions. *International Journal of Modern Physics*, **C12**:759–789, 2000.
- [29] M. Burgess, H. Haugerud, T. Reitan, and S. Straumsnes. Measuring host normality. *ACM Transactions on Computing Systems*, 20:125–160, 2001.
- [30] G. Box, G. Jenkins, and G. Reinsel. *Time Series Analysis*. Prentice Hall, New Jersey, 1994.

- [31] M. Burgess. Two dimensional time-series for anomaly detection and regulation in adaptive systems. *IFIP/IEEE 13th International Workshop on Distributed Systems: Operations and Management (DSOM 2002)*, page 169, 2002.
- [32] W.D. McComb. *Renormalization Methods: A Guide for Beginners*. Oxford University Press, 2003.
- [33] J. Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann, San Francisco, 1988.
- [34] M. Steinder and A. Sethi. A survey of fault localization techniques in computer networks. *Science of Computer Programming*, 53:165, 2003.
- [35] M. Burgess. Cfengine www site. <http://www.iu.hio.no/cfengine>, 1993.
- [36] Snort. Intrusion detection system. <http://www.snort.org>.
- [37] R.O. Duda, P.E. Hart, and D.G. Stork. *Pattern classification*. Wiley Interscience, New York, 2001.
- [38] R. Durbin, S. Eddy and A. Krigh, and G. Mitcheson. *Biological Sequence Analysis*. Cambridge, Cambridge, 1998.
- [39] B.J. Oommen and R.L. Kashyap. A formal theory for optimal and information theoretic syntactic pattern recognition. *Pattern Recognition*, 31:1159, 1998.
- [40] H. Bunke and J. Csirik. Parametric string edit distance and its application to pattern recognition. *IEEE Transactions on Systems, Man and Cybernetics*, 25:202, 1995.
- [41] S-H. Han, M-S Kim, H-T Ju, and J.W-K. Hong. The architecture of ng-mon: A passive network monitoring system for high-speed ip networks. *IFIP/IEEE 13th International Workshop on Distributed Systems: Operations and Management (DSOM 2002)*, page 16, 2002.
- [42] P. A. Porras and P. G. Neumann. EMERALD: Event monitoring enabling responses to anomalous live disturbances. In *Proc. 20th NIST-NCSC National Information Systems Security Conference*, pages 353–365, 1997.
- [43] Peter G. Neumann and Phillip A. Porras. Experience with EMERALD to date. pages 73–80.
- [44] A. Somayaji and S. Forrest. Automated reponse using system-call delays. *Proceedings of the 9th USENIX Security Symposium*, page 185, 2000.
- [45] P.D’haeseleer, S. Forrest, and P. Helman. An immunological approach to change detection: algorithms, analysis, and implications. In *Proceedings of the 1996 IEEE Symposium on Computer Security and Privacy (1996)*.
- [46] P. D’haeseleer. An immunological approach to change detection: Theoretical results. In *9th IEEE Computer Security Foundations Workshop (1996)*.
- [47] K.M. Begnum and M. Burgess. A scaled, immunological approach to anomaly countermeasures. *Proceedings of the VIII IFIP/IEEE IM conference on network management*, 2003.

- [48] R. Sekar, T. Bowen, and M. Segal. On preventing intrusions by process behaviour monitoring. *Proceedings of the workshop on intrusion detection and network monitoring*, USENIX, 1999.
- [49] B.D. Ripley. *Pattern recognition and neural networks*. Cambridge University Press, Cambridge, 1996.
- [50] J.A. Freeman and D.M. Skapura. *Neural networks: algorithms, applications and programming techniques*. Addison Wesley, Reading, 1991.
- [51] T. Lane. Machine learning techniques for the computer security. phd thesis, purdue university, 200.
- [52] M. Steinder and A. Sethi. Distributed fault localization in hierarchically routed networks. *IFIP/IEEE 13th International Workshop on Distributed Systems: Operations and Management (DSOM 2002)*, page 195, 2002.
- [53] P. Barford, J. Kline, D. Plonka, and A. Ron. A signal analysis of network traffic anomalies, 2002.
- [54] M. Sloman. Policy driven management for distributed systems. *Journal of Network and Systems Management*, **2**:333, 1994.
- [55] N. Damianou, N. Dulay, E.C. Lupu, and M. Sloman. Ponder: a language for specifying security and management policies for distributed systems. *Imperial College Research Report DoC 2000/1*, 2000.
- [56] H. Lewis and C. Papadimitriou. *Elements of the Theory of Computation, Second edition*. Prentice Hall, New York, 1997.
- [57] Ehab Al-Shaer and Hazem Hamed. Firewall policy advisor for anomaly detection and rule editing. In *Proc. IEEE/IFIP 8th Int. Symp. Integrated Network Management (IM 2003)*, pages 17–30, Mar 2003.