

Network Events in a Large Commercial Network: What can we learn?

Antoine Messenger*, George Parisi*, Robert Harper†, Philip Tee‡, István Z. Kiss* and Luc Berthouze*

*School of Engineering and Informatics, University of Sussex, Brighton, UK

Email: {a.messenger, g.paris, i.z.kiss, l.berthouze}@sussex.ac.uk

†Moogsoft Ltd, River Reach, 31-35 High Street, Kingston-Upon-Thames, UK, Email: rob@moogsoft.com

‡Moogsoft Inc, 1265 Battery St, San Francisco, CA 94111, Email: phil@moogsoft.com

Abstract—ISP and commercial networks are complex and thus difficult to characterise and manage. Network operators rely on a continuous flow of event log messages to identify and handle service outages. However, there is little published information about such events and how they are typically exploited. In this paper, we describe in as much detail as possible the event logs and network topology of a major commercial network. Through analysing the network topology, textual information of events and time of events, we highlight opportunities and challenges brought by such data. In particular, we suggest that the development of methods for inferring functional connectivity could unlock more of the informational value of event log messages and assist network management operators.

I. INTRODUCTION

Today’s commercial and ISP networks are large, complex, heterogeneous and fast-evolving. They usually contain a large amount of servers placed in data centres that are made of commodity, off-the-shelf network equipment to interconnect them and make them accessible to the world. They span a large number of geographical regions (commonly across continents) and provide end-users with access to network and content services and applications. Middleboxes, which implement a diverse set of in-network functionality, are also crucial for the provision of efficient and secure services [1]. Ensuring continuous service availability for such networks is extremely challenging [2]. Despite continuous innovation in the network data and control planes, innovation in the management plane has been slower [3].

Network operators rely on the continuous collection of event data from all network devices (including servers and workstations at the edges of the network) that are deemed to be important; this is most commonly a very large amount of devices. Event data is then collected and analysed either in real-time or off-line, by specialised software.

A key requirement of any such event analysis technology is the ability to identify (or even predict) an outage as quickly as possible before user experience gets disrupted. This, however, is rendered extremely challenging by two features of the data. First, network events are commonly produced at a very high rate, up to 10^6 events per second [4] making the outage identification process both time- and resource-intensive [5]. Second, the vast majority of these events are just noise and only a few of them correlate to ‘actionable events’.

A lot of research has been done on algorithms performing Root Cause Analysis (RCA) [6]; i.e. identifying network events that escalate to actual network and service outages. Steinder et al. identify three main approaches to extracting the causal sources of events. *Rule-based analysis* is commonly used [7]. In this approach, the network operator predefines a set of rules to identify the causal sources of events and exclude uninteresting devices. Rules need to be updated when there are changes in the network. Knowledge from the network topology and temporal correlations between events is exploited to make those approaches more accurate and flexible [8]. In *model traversing techniques* one explores progressively the neighbours of each entity emitting an event to identify its source [9] using a formal representation of the network structure. Finally, *graph-theoretic approaches* have also been proposed. They require a priori knowledge of how failure of one device affects other devices in order to build a causality or dependency graph. Those graphs are then used to return a number of fault hypotheses that best explain the observations. This has been shown to be an NP-hard problem and several heuristics have been developed to reduce the complexity of the problem [10], [11], [12]. Most of those approaches require a priori knowledge about the network and are therefore static. They are ill-suited to networks with multiple layers of logical connectivity and could be greatly improved by the use of temporal correlation between events [13]. They scale poorly with the number of events to be processed, therefore approaches to eliminate insignificant events have been proposed, by manually setting filters that exclude specific devices or types of events from processing by an RCA algorithm. Such static approaches are problematic in dynamic and fast-evolving networks like the ones of modern service and content providers. Recently, dynamic approaches for filtering network devices based on the notion of graph vertex entropy [14], [4] and supervised machine learning [15] have been proposed.

Despite this literature on efficient network event processing, there is actually little published information on the characteristics of network events from both in-network and end-host devices in large, modern, commercial networks. Heterogeneity in the functionality, and therefore the type of network events produced, makes RCA even more challenging. In this paper, we present a study of network events collected from a large commercial network over a period of five months. Through

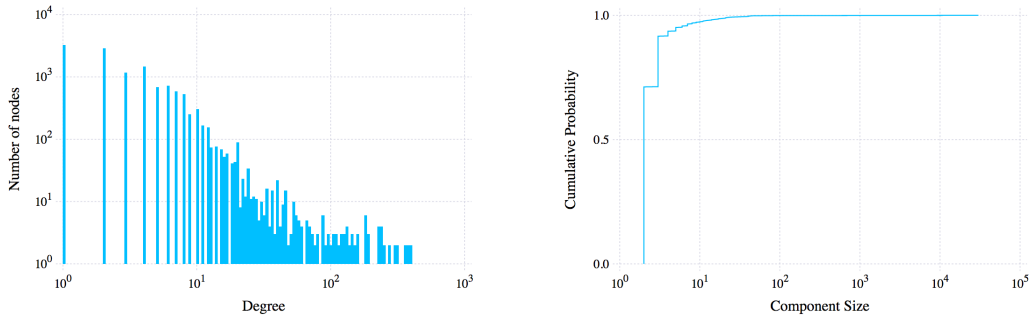


Fig. 1. Network topology: (left) node degree distribution for the largest connected component and (right) component size distribution.

attempting to reconcile information gathered from both event-related data and the underlying network topology, we aim to complement the community’s knowledge about this type of data as well as illustrate some of the challenges faced by network management operators with the view to stimulate new research in that area. In particular, we will argue that network management operators may benefit from methods that can use event logs to infer *functional connectivity*, i.e., a logical topology involved in the provision of a particular service, irrespective of the underlying physical topology. For example, a cluster of load balancers, application and database servers would comprise a meshed functional connectivity, where all nodes are logically connected to one another. We use this dataset to identify requirements and challenges for such methods.

II. DATASET DESCRIPTION

This paper is based on a dataset of network events and the underlying network topology, from a large commercial network¹. The network is comprised of a core of meshed backbone routers and a distribution layer of switches. It is primarily used to support the commercial operations of a Fortune 500 technology company, including accounting, human resources, research and development, communications and telephony and sales support activities. There are points of access to this network at almost every major city in the United States and multiple cities in most European and Asian countries. Network events span a duration of five months (2/5/2015 to 25/11/2015). The network topology was recorded at the end of that period.

Event log messages are collected from various sources in the network (end-hosts, switches, routers and middleboxes) and refer to functionality at various network layers; e.g. application-layer notifications, such as runtime exceptions, as well as network and link layer, such as routing protocol errors and links going up or down. *Emitted events* that are recognised as pertaining to the same network issue, warning or notification are then grouped into a so-called network *alert*. An alert is a Moogsoft construct² defined as “a de-duplicated event or

an instance of new data”. Our dataset contains alerts, which consist of a human-readable type and description, the number of de-duplicated events contained in the alert and timestamps for when the first and last events for the alert were collected. An alert is the user-presentable notion of a potential problem and is therefore handled as a ticket that needs to be ‘actioned’ by an operator. Events can be extracted from alerts although in this dataset not all event times will be known.

III. NETWORK TOPOLOGY

Network management providers may not always be in a position to maintain an accurate view of the network topology of their clients. Indeed, although a client’s devices and software (along with in-house management systems) are configured to emit events to the provider’s server(s), changes in the underlying topology are not necessarily visible to the network management provider. This is typically because accurate knowledge of the network topology is not necessary to run the operations’ monitoring software. Additionally, it may also be challenging to acquire the network topology, especially when changes in it are frequent. Finally, providing access to such information presents significant security risks. However, we will argue in Section VI that it could be automatically inferred from the events emitted by network devices.

The topology we were provided with was obtained by taking a dump of the customer’s operations database. This database is fed by the change management and connectivity discovery systems. It is automatically updated when network links are provisioned/de-provisioned or equipment is configured. The network consists of 73,677 devices and 117,846 physical connections among these devices. The underlying graph is not fully connected; its largest connected component (referred to as giant component thereafter) consists of 30,229 devices, which is less than half of the total number of devices. Figure 1 shows the degree distribution of the giant component (left panel) as well as the size distribution of the connected components (right panel). In the latter, we observe a very large number of very small components ($\sim 15,000$ components of size less than 10 devices each). This is a striking observation which, at first sight, could suggest a seriously flawed topology collection and recording process. However, in Section IV-A, we will show how an analysis of the recorded events enables

¹The dataset is currently not publicly available due to its commercially sensitive nature.

²<https://docs.moogsoft.com/display/060000/Glossary>

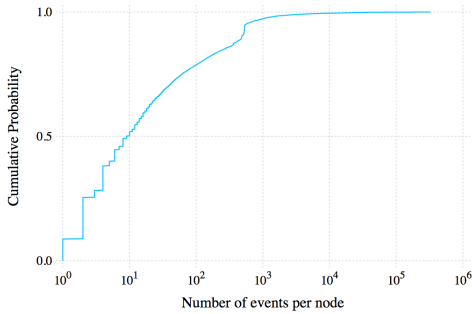


Fig. 2. Distribution of the number of emitted events per node

the formulation of hypotheses for such a large number of disconnected components. Importantly, an analysis of the giant component alone reveals key network metrics that are consistent with those typically found in large communication networks [16], namely: scale free-like nature of its degree distribution (under the “super weak” definition of scale-freeness proposed by Clauset et al. [17]), disassortativity ($r = -0.22$) and almost-zero clustering coefficient ($c = 0.00011$).

We conclude this Section with the following two observations: (a) a large fraction of nodes ($\sim 67\%$) in the recorded topology do not have any events associated with them; (b) a large number of events are associated with nodes that are not listed in the recorded topology. With respect to the former observation, we speculate that these are ‘uninteresting’ devices (probably workstations at the edges of the network), which are not configured to communicate with the event processing server but have been recorded in the topology nevertheless. On the other hand, the latter observation highlights the incomplete nature of the topology record, with a part of the network, including a large number of edge devices, not recorded despite its devices being configured to emit events to the server. However, this could also be because nodes that emitted events were no longer in the network (or their identifier changed) when the network topology was recorded. Note that events from these devices are fully usable when it comes to identifying faults with network components or applications, highlighting the fact that partial only knowledge of the topology does not hinder the operation of the network management system. In Section VI we will suggest that such events may provide the substrate to infer missing nodes as well as uncover information about changes in the network topology over time.

IV. NETWORK EVENTS AND ALERTS

Among the 53,604 network devices that emit events, 23,919 are part of the recorded network topology ($\sim 32\%$ of the total number of nodes in the recorded network topology) and 13,025 are part of the giant component ($\sim 42\%$ of the total number of nodes in the giant component). As mentioned in the previous section (and indicated by the percentages above), a large number of nodes in the recorded topology do not emit any events. We considered 21,223,756 events grouped in 473,580 alerts. An analysis of the distribution of number of

emitted events per node, shown in Figure 2, reveals that most nodes only emit a few events (~ 2 events per day). However, there is heterogeneity in behaviour and a few devices in the network emit a large number of events. The average number of emitted events per node indicates that in this dataset there is sparsity of information at node level. This sparsity is not particularly consistent with published numbers but is specific to the particular network presented here. More importantly, it has computational and mathematical implications which we will discuss in Sections V and VI.

A. Event Types

Every emitted event is accompanied by a type field which usually contains a single word. This field is mapped to specific types of events at deployment time. Such mapping can be manual, e.g. for events received from third-party event management/monitoring systems, or automatic, i.e., involving text processing to extract meaningful types from raw log messages. Out of a total of 41 distinct types³, 17 are found to account for over 98% of all events in the record. These are shown in Figure 3(left), with type “other” aggregating the remaining 24. For each type, we plot the number of emitted events, recorded alerts and the number of network devices involved in those.

The types of emitted events appear to fall under two categories; those that refer to in-network functionality (‘LAN-Switch’, ‘Router’) and those that refer to end-host/server functionality (all other types). The former account for almost half of all events (41%) whilst the three most prevalent types in the latter category (‘JVM’, ‘NT’ and ‘Linux’) account for 36% of the events.

Devices emitting ‘LANSwitch’ or ‘Router’ events account for only 6.6% of the total number of devices. This is consistent with the assumption that those devices are in-network devices (e.g., routers and switches) as in such a network we can expect the number of devices at the edges of the network to be significantly larger than that of in-network devices. Using the provided network topology, and subject to the caveat of nodes emitting events not being recorded, we found that there is very little overlap between sets of nodes involved in the emission of events of a particular type. Specifically, only 0.5% of all devices emit events of more than one type. In the main, devices seem to be emitting events that pertain to their functionality in the network. For example, we identified that events of type JVM are emitted by web back-end systems when failing to connect to other back-end servers (e.g. database servers). It is possible, of course, that the reason why devices only emit events of a single type can be attributed to choices made by the network management operator. In addition to having a potential impact on the event rate (i.e., this deployment may be less verbose than others), this consideration also has implications for the notion of functional connectivity as we will discuss in Section VI.

Devices within the giant component emit events of 13 (out of all 41) types, namely: ‘Linux’ (8074) - ‘VMWare’ (2833)

³Roughly 0.5% of the emitted events have an unspecified type field and were therefore excluded from further analysis.

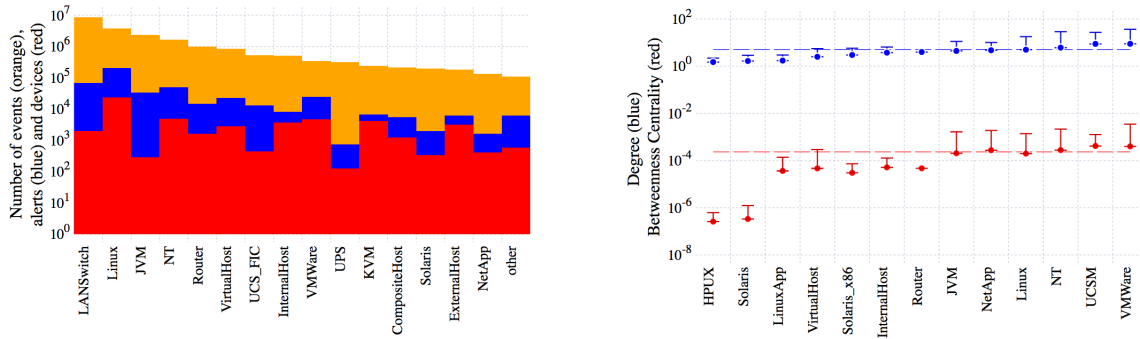


Fig. 3. (Left) Number of devices (red), recorded alerts (blue) and emitted events (orange) for each event type. (Right) Average degree and betweenness centrality (and their upper standard deviation) of nodes of the giant component who emitted events, organised by event type. For comparison purposes, the dashed lines show the mean degree and betweenness centrality for the full giant component (i.e., including devices that did not emit any event). Lower standard deviations are not shown due to log-scale plotting.

- ‘NT’ (862) - ‘VirtualHost’ (372) - ‘JVM’ (289) - ‘NetApp’ (228) - ‘UCSM’ (108) - ‘Solaris’ (15) - ‘InternalHost’ (12) - ‘LinuxApp’ (11) - ‘HPUX’ (2) - ‘Solaris_x86’ (2) - ‘Router’ (1). A striking observation is that none of these types pertain to in-network functionality (with the exception of the single event of type ‘Router’). This observation is counter-intuitive at first. However, juxtaposed with our previous observation in Section III that the recorded topology consisted of a large number of small disconnected components, it begs the question of whether some of the devices in the disconnected components may be duplicate entries of devices that are in the giant component. This is indeed plausible given that (a) routers (and their connectivity with other network devices) are typically known by, and recorded with, multiple IP addresses (and respective canonical names) that correspond to interfaces to different IP subnets they interconnect; (b) network switches are also known by, and recorded with, multiple IP addresses for administrative and monitoring purposes; and (c) it is common to use different IP addresses for network management functionality, i.e., for running SNMP and communicating with custom network management software. In Section VI, we will appeal to the network science notion of *multiplex networks* and argue that a suitable functional connectivity inference method should make it possible to match these devices.

Within the giant component itself, we find that devices that emit alerts show great variation in their degree and betweenness centrality (see right panel of Figure 3), suggesting that devices associated with a type are not randomly distributed within the network. For instance, devices that emit ‘VMWare’-type events tend to be centrally located and more connected than devices emitting events of other types. This yet again brings to fore the notion of functional connectivity.

B. Alert Descriptions

The description field of each event contains textual information produced by the source of the event (e.g., a specific Java exception that is pushed to the management system through the *stderr* stream). The vocabulary in the descriptions is very context-specific, containing specialised words and technical

jargon (e.g. “host or interface down”, “Database Instance Session exceeding temp threshold”). We used a custom iterative filtering approach to produce a minimal classification of the descriptions. At each iteration, we visually inspected a randomly selected set of unclassified alerts. From their description, we constructed a meaningful pattern of words (within the specific context) as basis for a new class. We automatically searched for all remaining alerts whose description matched the pattern and assigned them to the newly created class. The process was continued until the chosen pattern could not classify more than 1% of the recorded alerts. This iterative process led us to identify 14 classes that account for 94% of the recorded alerts (or 97% of the emitted events). Ninety six percent of all devices are involved in at least one of these classes, which are as follows:

- 1) Interface:= [name] down (e.g. interface:=eth1 down)
- 2) Host down: [canonical name]
- 3) Threshold Notification (e.g. [cpu10]=< 90.00)
- 4) Application down: [name] (e.g. application down: NTP)
- 5) Host CPU usage (e.g. Virtual Machine/Host CPU usage)
- 6) Server Inaccessible (along with chassis and blade information)
- 7) Auto cleared (e.g. AutoCleared due to inactivity timeout: null)
- 8) Memory usage problem
- 9) Database problem
- 10) Oracle related problem
- 11) Host not connecting to VC (e.g. ESX host [name] may not be connecting to VC [name])
- 12) Host unreachable (e.g. [name] Agent is unable to communicate with the OMS (agent is unreachable))
- 13) User logging out (e.g. User [name] logged out [time])
- 14) Device rebooting (e.g. [name] rebooted [time])

As shown by the left panel of Figure 4, classes ‘Interface Down’ and ‘Host Down’ account for 24% of all recorded alerts (50% of all emitted events). Furthermore, these two classes, along with classes ‘Threshold Notification’, ‘Oracle Problem’ and ‘Application Down’, which are notifications or describe

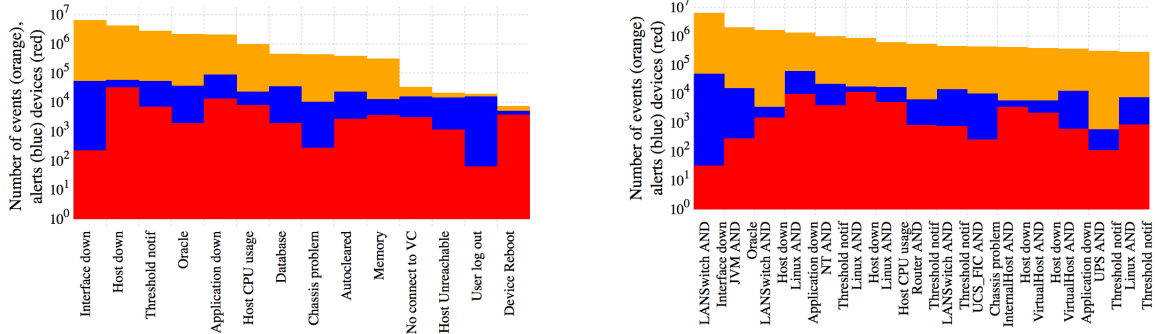


Fig. 4. Number of devices (red), alerts (blue) and events (orange) for each class of alert (left) and for each of the 15 most populated type-class pairings (right).

application layer events, account for 38% of the recorded alerts (87% of the emitted events). From a number of devices perspective, the ‘Host Down’ class is the most common, with 59% of all devices emitting events of this type. On the other hand, there are only 228 devices that emit events of class ‘Interface Down’ class.

An analysis of all type-class pairings revealed that just 15 pairings (right panel of Figure 4) account for 53% of the total number of recorded alerts (80% of the total number of events). Strikingly, one of those (LANSwitch and Interface Down) account for 10% of all alerts (33% of all events) despite only involving 33 devices, none of which are located in the giant component (we refer the reader to Section IV-A for potential reasons why devices that emit events relevant to in-network functionality are not part of the giant component). Events of this specific pair indicate links that go down, as identified by neighbouring switches, and are unsurprisingly very common, given the large size of the network, including end-devices.

V. ANALYSIS OF TEMPORAL CORRELATIONS

Since the presence of temporal correlations between emitted event times can provide insights into possible functional relationships between the devices that emit them [6], we used standard correlation techniques (see [18], for example) to gain further insights into the dataset. Namely, we assessed the presence of a functional relationship between two devices in terms of whether the Pearson’s correlation coefficient between the time series of their respective recorded alert timestamps significantly differed from that expected under the null hypothesis that they were independent. Formally, this process involves calculating the Pearson’s correlation coefficient and applying the Fisher’s z -transformation. A pair of time series of recorded alerts was considered to be significantly correlated if their z value was greater than one standard deviation (95% confidence level) [19].

In this dataset only recorded alerts were provided with timestamps and therefore our analysis relied on alerts rather than events. Further, since most network devices are linked to a single type of recorded alerts, we associated each device to its most frequent alert type. Note that due to the very low node-

level alert rate reported earlier, binning was needed before analysis. An arbitrary bin size of one hour was used.

Since little could be expected from a pairwise analysis due to the extreme sparsity of the data, we only report summary statistics based on a group-level analysis. In what follows, we will refer to as *within-type pair* a pair of devices that have the same most frequent type of recorded alerts, and *between-type pair*, a pair of devices whose most frequent types of recorded alerts are distinct. For each of the alert types, we calculated the proportion of within-type pairs who showed significant correlations in the time series of their respective recorded alert timestamps and compared it to the proportions of between-type pairs that showed significant correlations. Statistical differences between within-pair proportion and between-pair proportions were assessed using a simple boxplot-based outlier detection method [20].

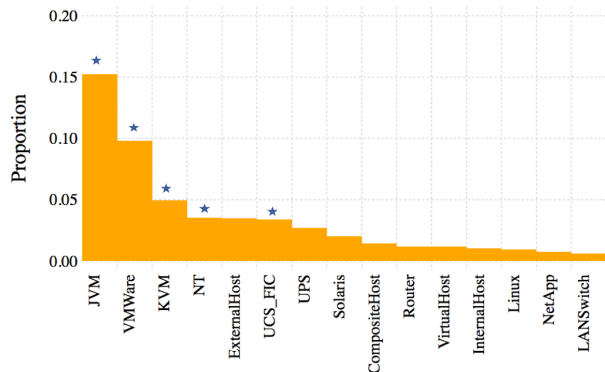


Fig. 5. Proportion of pairs of devices whose alert types are identical and whose alert times are significantly correlated, for each of the 15 most common types. An asterisk denotes the fact that this proportion is significantly higher than expected when considering all possible pairs of devices, irrespective of type.

This analysis revealed that for 18 of the 41 types, within-type pairs, i.e., pairs of devices with recorded alerts of the same type, were more likely to have correlated alert times than between-type pairs, i.e., pairs of devices with recorded alerts of a different type. On the one hand, this increased

likelihood is not particularly surprising (especially given that in this dataset devices tend to only emit events of a single type) and provides support to the notion that temporal correlations could enable the identification of devices belonging to the same functional topology. On the other hand, the analysis also strongly suggests that these correlations are not trivial. First, not all event types feature significant correlations in the alert times of their devices. For example, types ‘Router’ and ‘LANSwitch’ (which account for 41% of all events) do not show significant correlations. Further, the 18 types that do show significant correlations only account for 29% of the alerts (25% of the events). Second, as shown by Figure 5, in those types that show significant correlations (see asterisk), less than 20% of the pairs show significant correlations. This therefore suggests that correlations are not trivial but may be localised spatially and temporally (see Section VI).

VI. DISCUSSION

Our analysis of a major commercial network has elicited a number of insights which we believe highlight the need for methods able to infer functional topologies within a network deployment based on the temporal characteristics of events emitted by its network devices.

Although the total number of events to be processed by a network management server can be really large, device-level event rates can be very low, which makes device-based analysis very challenging. The reasons for such sparsity are multiple and include whether the device exists at all times during the record, whether or not it is monitored for specific types of events, and indeed whether it is affected by such types of event. In this deployment, network devices tended to emit events of a single type. Although this increased sparsity, it also helped with identifying the role of the devices in the network (e.g., application or database servers, routers or switches) which can be crucial to building a reliable view of the underlying physical topology.

Getting access to reliable information about the physical topology is not straightforward. Yet, such knowledge would be beneficial for a network management provider that would otherwise have to rely solely on incoming events to identify service outages in an unknown (or partially known) network topology. The presence of disconnected components presents a challenge for a traditional graph-based analysis. An alternative approach is to think of the network as a *multiplex network*, i.e., a network where the same set of nodes are linked by different types of interaction. In this approach, each disconnected component may provide a graph-representation of a given layer of interaction. What is needed then is a means to match the nodes across layers of interaction.

Our analysis of time series of recorded alerts in Section V hints at a possible approach to elicit layers of interaction, by showing that correlations exist for alerts of the same type, recorded for different network devices. Although a much more robust framework is needed (e.g., to deal with sparsity and time-varying nature of the signals), the existence of correlations opens up the possibility of inferring

common functionality of different devices in the network, i.e., functional topologies. Such functional topologies may be application specific (e.g. a web application development that consists of application and database servers), refer to in-network functionality (e.g. an OSPF area), or cut across layers (e.g. when switches report failed links to a database server while application servers try to connect to the same database server). This would provide an extremely powerful framework for network management providers to identify (or even predict) service outages, since functional topologies can be thought of as defining spheres of influence for network devices whereby the operational status or (mis)behaviour of one or more devices would influence the operational status or behaviour of other devices within its sphere of influence.

REFERENCES

- [1] R. Potharaju and N. Jain, “Demystifying the Dark Side of the Middle: A Field Study of Middlebox Failures in Datacenters,” in *Proc. of ACM IMC*, 2013.
- [2] R. Govindan, I. Minei, M. Kallahalla, B. Koley, and A. Vahdat, “Evolve or Die: High-Availability Design Principles Drawn from Googles Network Infrastructure,” in *Proc. of ACM SIGCOMM*, 2016.
- [3] A. Gember-Jacobson, W. Wu, X. Li, A. Akella, and R. Mahajan, “Management Plane Analytics,” in *Proc. of ACM IMC*, 2015.
- [4] P. Tee, G. Parisi, and I. Wakeman, “Vertex Entropy As a Critical Node Measure in Network Monitoring,” *IEEE Transactions on Network and Service Management*, vol. 14, no. 3, pp. 646–660, 2017.
- [5] S. Kobayashi, K. Fukuda, and H. Esaki, “Mining causes of network events in log data with causal inference,” in *Proc. of IFIP/IEEE IM*, 2017.
- [6] M. Steinder and A. S. Sethi, “A survey of fault localization techniques in computer networks,” *Science of Computer Programming*, vol. 53, no. 2, pp. 165–194, 2004.
- [7] A. Patel, G. McDermott, and C. Mulvihill, “Integrating network management and artificial intelligence,” in *Proc. of IFIP/IEEE IM*, 1989.
- [8] K. Appleby, G. Goldszmidt, and M. Steinder, “Yemanja-a layered event correlation engine for multi-domain server farms,” in *Proc. of IFIP/IEEE IM*, 2001.
- [9] S. Kätker and M. Paterok, “Fault isolation and event correlation for integrated fault management,” in *Proc. of IFIP/IEEE IM*, 1997.
- [10] A. T. Bouloutas, S. Calo, and A. Finkel, “Alarm correlation and fault identification in communication networks,” *IEEE Transactions on Communications*, vol. 42, no. 234, pp. 523–533, 1994.
- [11] I. Katzela and M. Schwartz, “Schemes for fault identification in communication networks,” *IEEE/ACM Transactions on Networking*, vol. 3, no. 6, pp. 753–764, 1995.
- [12] S. Kliger, S. Yemini, Y. Yemini, D. Ohsie, and S. Stolfo, “A coding approach to event correlation,” in *Proc. of IFIP/IEEE IM*, 1995.
- [13] M. Hasan, B. Sugla, and R. Viswanathan, “A conceptual framework for network management event correlation and filtering systems,” in *Proc. of IFIP/IEEE IM*, 1999.
- [14] P. Tee, I. Wakeman, G. Parisi, J. Dawes, and I. Z. Kiss, “Constraints and entropy in a model of network evolution,” *The European Physical Journal B*, vol. 90, no. 11, p. 226, Nov 2017.
- [15] R. Harper and P. Tee, “The application of neural networks to predicting the root cause of service failures,” in *Proc. of IFIP/IEEE IM*, 2017.
- [16] M. Faloutsos, P. Faloutsos, and C. Faloutsos, “On power-law relationships of the internet topology,” in *Proc. of ACM SIGCOMM*, 1999.
- [17] A. Clauset, C. R. Shalizi, and M. E. J. Newman, “Power-law distributions in empirical data,” *SIAM Rev.*, vol. 51, no. 4, pp. 661–703, Nov. 2009.
- [18] A. A. Mahimkar, Z. Ge, A. Shaikh, J. Wang, J. Yates, Y. Zhang, and Q. Zhao, “Towards automated performance diagnosis in a large IPTV network,” in *ACM SIGCOMM Computer Communication Review*, vol. 39, no. 4, 2009, pp. 231–242.
- [19] S. K. Kachigan, *Statistical analysis: An interdisciplinary introduction to univariate & multivariate methods*. Radius Press, 1986.
- [20] J. W. Tukey, “Exploratory data analysis,” pp. 39–49, 1977.