# Event-Driven Ontology Population - From Research to Practice in Critical Infrastructure Systems[*]

David Graf[1,2], Wieland Schwinger[1], Werner Retschitzegger[1], Elisabeth Kapsammer[1], and Norbert Baumgartner[2]

[1] Johannes Kepler University, Linz, Austria
firstname.name@cis.jku.at
[2] team GmbH, Vienna, Austria
firstname.name@te-am.net

**Abstract.** In an interconnected world, the *Systems-of-Systems (SoS)* paradigm is prevalent in various domains, particularly in large-scale environments as being found in the area of critical infrastructures. Due to the characteristics of SoS and those of critical infrastructures, the realization of high level services for Operational Technology Monitoring (OTM), such as failure cause reasoning, is challenging, whereas *interoperability* and *evolvability* are most pressing. In this realm, the contribution of this paper is twofold: Firstly, we conduct a systematic literature review focusing on semantic technologies in areas like (i) semantic annotations, (ii) event log focused work in the IoT, (iii) organizational process mining, and (iv) complex event processing. Based thereupon, we elaborate towards a hybrid (semi)-automatic ontology population approach in the context of OTM by combining inductive and deductive methods.

**Keywords:** Systems-of-Systems, Critical Infrastructure Systems, Ontology Population, Operational Technology Monitoring

## 1 Introduction

*Operational Technology Monitoring (OTM).* In an interconnected world, the *Systems-of-Systems (SoS)* paradigm is prevalent in various domains, particularly in large-scale environments as being found in the area of critical infrastructures. One example of such a *Critical Infrastructure System (CIS)* are *Intelligent Transportation Systems (ITS)*. By the interplay of various constituent systems (e.g, a video system, electronic road signs, or tunnel operating programs), ITS's aim is to provide services on top (e.g., efficient monitoring and controlling of traffic) based on a wide range of technologies, aka. *Operational Technologies (OT)* being more and more based on Internet-of-Things (IoT) [25]. Due to the high demands on the reliability of these services, adequate techniques for monitoring

the wide range of OT used, i.e., *Operational Technology Monitoring (OTM)* are needed. In this respect, OTM is responsible to ensure high availability of OT, for instance, to reactively or even proactively trigger maintenance actions in case of an identified (prospective) failure of OT.

*Challenges in CIS.* The realization of high level services for OTM (e.g., service quality monitoring of OT [10] or failure cause reasoning) is, however, due to the characteristics of SoS [21] and those of critical infrastructures [6], challenging. Constituent systems, mostly geographically distributed, are focusing on very specialized and encapsulated tasks, often operating in an isolated manner, which lead to massive heterogeneity at different levels [27]. Hence, *interoperability* is lacking and thereby an integrated view of OTM across systems is very limited. This is aggravated by prevalent legacy systems, since building up the entire infrastructure from scratch is not the standard-case in CIS, predominately showing heterogeneity of data lacking structured and semantic information. In addition, the dynamic nature of such CIS environments lead to omnipresent evolution of systems comprising behavioral aspects (e.g., OT failures, OT maintenance) as well as structural aspects, meaning that the underlying OT is continuously added, removed or replaced. Thus, there is the need to deal with *evolvability* within OTM, which dramatically increases its complexity. These key challenges hamper an integrated and up-to-date view, i.e., a conceptual representation, of the entire OT infrastructure, being independent from the actual technology used, which is, however, an indispensable prerequisite for enabling efficient OTM and providing high level services for OTM on top.

*(Semi-)Automatic Ontology Population.* A promising paradigm to address these challenges are *semantic technologies* in terms of *ontologies* [9]. While the ontology's *T-Box* can be manually specified by domain experts through modeling the OT objects and the relationships in between at type-level, it is simply not feasible from a practical point of view to manually populate an ontology's *A-Box* with hundred thousands of objects and their links in between. To give an example from practice, the national highway network focused by our work and being the underlying example throughout this paper comprises more than 100.000 OT devices of more than 200 different types, ranging from simple sensing and actuating devices (e.g., a CO-sensor) to more complex systems consisting of many devices of various types (e.g., a video system), which are geographically distributed over 2.220 kilometers highway and 165 tunnels (cf. our previous work, e.g. [11]). Thus, (semi-)automatic ontology population is a must, especially in the light of *evolvability* outlined above.

*Contribution and Paper Structure.* In order to enhance (semi-)automatic population of an OT ontology's A-Box, the current paper's contribution is twofold: In Section 2, we conduct a systematic literature review of most promising approaches of related research areas. Based thereupon, we elaborate towards a hybrid (semi)-automatic ontology population approach in the context of OTM by combining inductive and deductive methods in Section 3. While using implicit knowledge provided by OT event streams as a basis, we gain additional infor-

mation on top by representing explicit knowledge provided by domain experts. We conclude and discuss future work in Section 4.

## 2   Promising Lines of Research

Addressing the primary goal of our work namely (semi-)automatic population of an ontology's A-Box with OT objects and their links in between, related work can be found in various areas most relevant in those of (i) *semantic annotation* from text or semi-structured data, surveyed by [20], (ii) event log focused work in the IoT such as *event log mining* and *event analysis*, (iii) the area of *organizational process mining* using semi-structured log data as data source to extract knowledge about underlying resources, most promising mining approaches surveyed by [22], as well as (iv) *complex event processing* as a specific form of *stream processing* dealing with the streaming aspect in such environments, most promising approaches surveyed by [1]. The following literature review considering these areas and compares related work primarily based on the *source* data structure, the *techniques* applied, and the *target* data structure. Table 1 gives a summary along these comparison dimensions.

*Semantic Annotation.* The work of [14] and [28] use data-driven techniques such as clustering and semi-supervised classification to populate a domain ontology, the latter being closely related regarding the target data structure by populating an ontology's A-Box with resources in terms of web services. Both, however, use primarily unstructured text documents as data source and do not use semi-structured stream data originating from event logs. Closely related to our work is the approach of [4] populating an event ontology's A-Box for monitoring vineyards grounded on a heterogeneous IoT sensor network aiming to mine causality relationships between events occurred during the life cycle of a wine production. Although the data originates also from IoT sensors, the target ontology primarily focus on events and not on a representation of the underlying IoT objects. Related with respect to techniques used are approaches in the area of semantic annotation, i.e., the process to annotate entities in a given text with semantics [8] (e.g., using ontology classes), such as the work of [18][19], which, however, address primarily the input data source themselves (often in terms of text documents) as target data structure.

*Event Log Mining.* While the work considered so far focus rather on the target data structure, event log mining and event analysis focus on the source data structure, i.e., dealing with huge amount of event log data generated on a daily basis in various systems. Data mining techniques and tools are commonly used, such as the work of [3] applying natural language processing and information extraction techniques by using the tool GATE in order to semantically enrich event log data. Moreover, [31] use the tool LogClusterC to mine and discover line patterns, similar to event types, from semi-structured textual event logs. Both, however, focus on enrichment of the log itself, rather than using information provided by the log for other purposes such as failure reasoning or populating an ontology. Closely related to ours is the work of [12] transforming air quality

sensor data to a semantic-based representation in terms of a Ressource Description Framework (RDF) model. In addition, their presented "event and clustering analytic server", considered as a middleware, is based on the OpenIoT platform [16] being a well-known semantic-based framework in the IoT context considered in one of our previous work [9]. The created ontology, however, models the event data itself, rather than representing the underlying resources, i.e., the individual sensors. With respect to mining correlations from event logs, the approach followed by [30] is worth to mention, introducing an abstract concept of a "service-hyperlink" which represents dependencies of data services based on the correlation of events in the log. Their focus, however, lies on dependencies between services rather than between individual devices of OT.

*Organizational Process Mining.* Aiming to "derive the underlying organizational structure of a CPS" from event logs, most promising mining approaches are reviewed by [22]. Discussed techniques such as "metrics based on (possible) causality" focusing on temporal succession of activities, or "metrics based on joint cases" focusing on frequency and correlation of resources, being suitable to derive relationships between objects, seems to be promising for our work, thus being considered as widely related. This is also the case for variations of distance measures, e.g, those used by [26], and traditional clustering techniques applied to event logs, e.g., those used by [15] as well as, since "time is a key relation between pieces of information" [7], time-based approaches such as the organizational mining approach of [13]. With respect to the source data structure and the techniques used, closely related is the approach of [5] in terms of semantically annotating event log information in order to enhance the discovery of unknown dependencies, i.e., how activities, or events respectively, are connected or who performed the activity, i.e. the resources respectively. Although these parts of their work are similar to ours, their focus lies on enhancing medical decision making rather than representing underlying resources.

*Complex Event Processing.* With respect to the streaming aspect prevalent in the environment considered by our work, we can find related approaches in the area of stream processing and one of its specialization namely complex event processing (CEP). The event based pattern matching approach TPStream based on Allens interval-algebra [2] proposed by the work of [17] identifies temporal relationships between events, however, their work focus on identifying complex temporal patterns rather than derive additional knowledge from rather simple patterns. Closely related is the work of [7] proposing a new semantic model for real-time reasoning from sensor data. As in our work they transfer event data to a semantic based model in terms of RDF triples, however, they represent context information of events, rather than derive knowledge of underlying resources. Related regarding identifying relationships between OT devices is the approach followed by [29] in terms of combining event information with background knowledge, however with the aim to improve processing quality rather than deriving relationships. Worth mentioning with respect to pattern learning is the approach of [23] using rule-based machine learning to learn new patterns prevalent in event data, thus considered as widely related to our work.

As Table 1 shows, although approaches discussed so far are related to our work in some of the comparison dimensions, however, none of them is directly applicable to our requirements since none of them aim populating an ontology's A-Box as target data structure based on event log data originating from message streams as data source structure. Hence, in the following, we propose a hybrid approach for populating an OT ontology's A-Box inspired by approaches discussed so far.

**Table 1.** Related approaches

| Category | Subcat | Item | Ganino et al. [8] | Lin et al. [18] | Liu et al. [19] | Jayawardana et al. [14] | Reyes-Ortiz et al. [28] | Belkaroui et al. [4] | Amato et al. [3] | Zhuge et al. [31] | Hromic et al. [12] | Zhu et al. [30] | Matzner u. Scholta [22]a | Ni et al. [26] | Jin et al. [15] | Jafari et al. [13] | Detro et al. [5]b | Endler et al. [7] | Korber et al. [17] | Teymourian et al. [29] | Mehdiyev et al. [23] |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Semantic Annotation | | | | | | Event Log Mining | | | | Process Mining | | | | | CEP | | | |
| Source | Structure | unstructured | ✓ | | | ✓ | ✓ | | ✓ | | | | | | | | | | | | |
| | | semi-structured | | ✓ | | | | ✓ | ✓ | ✓ | | | | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ | |
| | | structured | | | | | | | | | ✓ | ✓ | | | | | | | | ✓ | ✓ |
| | Data | at-rest | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | | | | ✓ | ✓ | ✓ | ✓ | ✓ | | | | |
| | | in-motion | | ✓ | | | | | | | ✓ | ✓ | | | | | | ✓ | ✓ | ✓ | ✓ |
| Techniques | Information Retrieval | text-based | | | | ✓ | ✓ | ✓ | ✓ | | | | | | | | | | | | |
| | | NLP | | | | | | | ✓ | | | | | | | | | | | | |
| | | IE | | | | | | | ✓ | | | | | | | | | | | | |
| | Machine Learning | supervised | | | | ✓ | | ✓ | | | | | | | | | | | | | ✓ |
| | | semi-supervised | | | | ✓ | | | | | | | | | | | | | | | |
| | | unsupervised | | ✓ | | | | | | | ✓ | ✓ | ✓ | ✓ | ✓ | | | | | | |
| | Data Analysis | distance-based | | | | ✓ | | | | | | | ✓ | | | | | | | | |
| | | similarity-based | | | | | | | | | | | | ✓ | ✓ | | | | | | |
| | | temporal-based | | | | | | | | | | | | | | | | | ✓ | | |
| | | others | | | | | | | | | | ✓ | | | | | | | | | |
| | Other Methods | tool-usage (e.g, GATE) | ✓ | | | | | | ✓ | ✓ | | | | | | | ✓ | | | | |
| | | CEP | | | | | | | | | | | | | | | | ✓ | ✓ | ✓ | |
| | | specific data-mining | | | | | | | | ✓ | | ✓ | | | | ✓ | | | | | |
| | | stream processing | | | | | | | | | ✓ | ✓ | | | | | | | | ✓ | ✓ |
| Target | Content-Ontology | wine events | | | | | | | ✓ | | | | | | | | | | | | |
| | | health events | | | | | | | | | | | | | | ✓ | | | | | |
| | | law events | | | | ✓ | | | | | | | | | | | | | | | |
| | | IoT data | | ✓ | ✓ | | | | | ✓ | | | | | | | | | ✓ | | |
| | | log events | | | | | | | ✓ | | | | | | | | | | | | ✓ |
| | Resource-Ontology | OT | | | | | | | | | | | | | | | | | | | |
| | | roles | | | | | | | | | | | | | | | ✓ | | | | |
| | | web services | | | | | ✓ | | | | | | | | | | | | | | |
| | Other Formalisms | org. model | | | | | | | | | | | ✓ | ✓ | ✓ | | | | | | |
| | | documents | ✓ | | | | | | | | | | | | | | | | | | |
| | | RBAC model | | | | | | | | | | | | | | ✓ | | | | | |
| | | event patterns | | | | | | | ✓ | | | | | | | | | | | ✓ | |

a survey, b preliminary work and not yet validated

# 3   Hybrid Ontology Population Approach

*Approach at a Glance.* Our approach is designed for real-time use, i.e., to be applied to a stream of human interpretable service-messages (e.g, a sensor sends a $CO_2$ value) and status-messages (e.g., a device notifies an error) consisting of (i) human interpretable message text, (ii) a unique identification of the affected OT device, (iii) the OT device type, and (iv) temporal information. The rational behind using human interpretable messages is that those messages can be seen

as the lowest common denominator of various heterogeneous systems, since all of them somehow communicate and interact with a human operator, i.e., send messages and receive control actions. The main idea of our hybrid approach is to combine inductive and deductive methods in order to populate an OT ontology's A-Box. Thereby, the inductive part is based on the event-information of logs, i.e., the messages, to (i) instantiate OT devices as objects, and (ii) compute the correlation between those devices allowing to (semi-)automatically derive potentially existing links among them. This inductive part is complemented by the deductive part of our approach in terms of gaining additional A-Box information based on the conceptualization of the domain, i.e., the T-Box (OT-types and their relationships in between), modeled by domain experts.

*Rational of our Approach.* The rational behind our proposed approach, or respectively, the reasons why a semantically rich OT ontology A-Box is beneficial for OTM purposes, is the following. Firstly, reasoning at A-Box level provides more information with respect to the monitored environment than reasoning on T-Box level, only, since the A-Box can be considered as a concretization of the T-Box (e.g., location information of individual OT objects showing spatial proximity to other OT objects). Secondly, meta-data about individual OT objects can be used for reasoning, which is crucial for OTM (e.g., hours of operation impacting the probability of a certain failure category). Thirdly, the OT ontology A-Box can be linked with context-information highly relevant for OTM (e.g., current weather or traffic situation impacting the criticality of an OT failure). Fourthly, links between objects allow identifying concrete other objects being affected by a failure or being the potential cause of a failure (e.g., the camera failure is caused by a temporarily overloaded server). In the following, we discuss the details of our approach based on the three core phases (i) correlation analysis, (ii) link instantiation, and (iii) object instantiation (cf. Fig. 1).

*Correlation Analysis (1).* In order to identify the correlation between OT devices (which further allow to derive links in the link instantiation phase), we adhere on techniques for calculating the correlation between individual OT devices based on temporal information of events (adapting and extending most promising work discussed above such as [17] to our specific purposes). The rational behind using techniques of (temporal) correlation analysis to derive links is that events of OT devices, which have some kind of correlation in-between to fulfill certain services, potentially occur in a certain degree of simultaneity or even in recurring patterns (e.g., a failure of an energy supply device will cause corresponding timely events of those devices being connected to that energy supply). Calculating the correlation (we adhere on similarity calculations based on "functional connectivity" of [24]) between the OT devices results in a n-times-n sized device-correlation-matrix (cf. Fig. 1), where n is the number of OT devices, showing the correlation between all possible devices. At this point, the inductive part of our approach gets in touch with the deductive part in terms of a plausibility-check, i.e., restricting the correlation calculations between devices where their corresponding types have a relationship modeled in the T-box, only. The device-correlation-matrix is continuously updated (considering a stream of

**Fig. 1.** Hybrid ontology population approach

messages) since new messages potentially provide new information and might change the correlation-values (e.g., the correlation-value between two OT devices increases if their corresponding message behavior show similar patterns, for instance, "cam_1" always sends a message before "srv_a").

*Link Instantiation (2).* During link instantiation phase the device-correlation-matrix is used as a basis to instantiate links between two OT objects in the A-Box. A link is instantiated and annotated with a probability, only, if the correlation value is above a certain threshold since otherwise links are too unstable (links are instantiated and deleted frequently when correlation values slightly changes). Additionally, the link must be consistent with the corresponding relationship modeled in the T-Box (e.g., given multiplicities). Since our approach is designed for a stream of messages, already existing links in the A-Box have to be recurringly rechecked and possibly updated (or even deleted) based on the correlation values (cf. center of Fig. 1). Thereby, link instantiation works in a loop-wise way, i.e., focusing the OT object, each relationship in the T-Box, or link in the A-Box respectively, is processed (instantiated, updated, or deleted) one after another (cf. examples visualized in Fig. 1).

*Object Instantiation (3).* While OT devices occurring in the message log can be instantiated straight forwardly to OT objects, existing but not yet populated OT devices (e.g. due to having not (yet) sent a message - aka "silent objects" and consequently "silent links") can be inferred through explicit knowledge provided by domain experts. The main idea is to consider T-Box information in a way that mandatory objects and their links are instantiated in the OT ontology's A-Box although not included in the message logs (e.g., based on T-Box information, a camera device is mandatory connected to a media server and therefore we are able to instantiate the "silent" media server object as well as the "silent link" to the camera object at the moment when the camera object is instantiated - cf. examples of Fig. 1). This method gives a series of benefits. Firstly, for monitoring purposes it enables to bridge „blind spots" in the monitored environment being the consequence of lacking integration of systems or even of unmonitored OT areas. Secondly, it enables identifying objects being affected by a failure or being the potential cause of a failure although those objects do not exist in the log (so far). Thirdly, in a semantically richer A-Box, object and link meta-data of silent objects can be used for reasoning purposes, also. Fourthly, instantiating silent objects and silent links works towards completing the A-Box. On the downside, this however comes at the costs that new objects derived from the event-information of the message-stream have to be merged with existing silent objects in the A-Box if the silent object was already instantiated beforehand. At this point the device-correlation-matrix again is an indication whether a new object is "the same" as a silent object and as a consequence have to be merged. Secondly, by considering T-Box relationship information to derive links, multiplicities have strong impact (e.g., T-Box relationships of *optional* multiplicity (0..1 or 0..*) can not be considered to derive "silent objects and silent links").

At this point we have to emphasize that the approach discussed so far is still work-in-progress meaning that we have some parts already implemented as

a prototype in a CIS domain adapted setting such as the correlation-analysis phase and the initial population of objects, whereas implementation of parts of the link instantiation phase is still ongoing. Nevertheless, first experimentations with real-world data containing 822k events (status- and service messages) from the year 2019 from a certain part of a highway network show promising results.

## 4    Conclusion and Future Work

The ontology population approach presented in this paper is grounded on a systematic literature review of related research areas dealing with similar requirements (especially regarding source and target data structure), which are (i) semantic annotations, (ii) event log focused work in the IoT, (iii) organizational process mining, and (iv) complex event processing. Based thereupon, we elaborate towards a hybrid (semi)-automatic ontology population approach in the context of OTM by combining inductive and deductive methods. Since the presented approach is work-in-progress, beside ongoing work at the core phases such as experimenting with configurations and parameters (e.g., correlation calculations or threshold values) as well as different techniques for A-Box updating strategies (e.g, interval-based methods), future work includes dealing with algorithm performance and computational complexity since the large amount of objects as well as the (dependent on the OT-type) high message frequency lead to high requirements on computation in real-world settings.

## References

1. Alevizos, E. et al.: Probabilistic Complex Event Recognition: a Survey. ACM Computing Surveys (CSUR) 50(5), 1–31 (2017)
2. Allen, James F.: Maintaining Knowledge about Temporal Intervals. Communications of the ACM 26(11), 832–843 (1983)
3. Amato, F. et al.: Detect and Correlate Information System Events through Verbose Logging Messages Analysis. Computing 101(7), 819–830 (2019)
4. Belkaroui, R. et al.: Towards Events Ontology Based on Data Sensors Network for Viticulture Domain. In: Proc. Int. Conf. on the IoT, pp. 1–7. ACM (2018)
5. Detro, S. et al.: Enhancing Semantic Interoperability in Healthcare Using Semantic Process Mining. In: Proc. Int. Conf. on Information Society and Technology, pp. 80–85 (2016)
6. Ellinas, G. et al.: Critical Infrastructure Systems: Basic Principles of Monitoring, Control, and Security. In: Intelligent Monitoring, Control, and Security of CIS, pp. 1–30. Springer (2015)
7. Endler, M. et al.: Towards Stream-based Reasoning and Machine Learning for IoT Applications. In: Intelligent System Conf., pp. 202–209. IEEE (2017)
8. Ganino, G. et al.,: Ontology Population for Open-Source Intelligence: A GATE-based Solution. Software Practice and Experience 48(12), 2302–2330 (2018)
9. Graf, D. and Schwinger, W. and Kapsammer, E. and Retschitzegger, W. and Baumgartner, N.: Cutting a Path Through the IoT Ontology Jungle - a Meta Survey. In: Int. Conf. on Internet of Things and Intelligence Systems. IEEE (2019)

10. Graf, D. and Schwinger, W. and Kapsammer, E. and Retschitzegger, W. and Pröll, B. and Baumgartner, N.: Towards Operational Technology Monitoring in ITS. In: Int. Conf. on Management of Digital Eco-Systems. ACM (2019)
11. Graf, D. and Schwinger, W. and Kapsammer, E. and Retschitzegger, W. and Pröll, B. and Baumgartner, N.: Towards Message-Driven Ontology Population - Facing Challenges in Real-World IoT. In: World Conf. on Information Systems and Technologies, pp. 361–368. Springer (2020)
12. Hromic, Hugo et al.: Real Time Analysis of Sensor Data for the IoT by means of Clustering and Event Processing. In: Proc. Int. Conf. on Communications, pp. 685–691. IEEE (2015)
13. Jafari, M. et al.: Role Mining in Access History Logs. J. of Computer Information Systems and Industrial Management Applications 1 (2009)
14. Jayawardana, V. et al.: Semi-Supervised Instance Population of an Ontology using Word Vector Embeddings. In: Proc. Int. Conf. on Advances in ICT for Emerging Regions, pp. 217–223. IEEE (2017)
15. Jin, T. et al.: Organizational Modeling from Event Logs. In: Proc. Int. Conf. on Grid and Cooperative Computing, pp. 670–675. IEEE (2007)
16. J. Kim and J. Lee: OpenIoT: An Open Service Framework for the Internet of Things. In: Proc. of World Forum on IoT (WF-IoT), pp. 89–93 (2014)
17. Körber, M. et al.: TPStream: Low-latency and High-throughput Temporal Pattern Matching on Event Streams. Distributed and Parallel Databases pp. 1–52 (2019)
18. Lin, S. et al.: Dynamic Data Driven-based Automatic Clustering and Semantic Annotation for IoT Sensor Data. Sensors and Materials 31(6), 1789–1801 (2019)
19. Liu, F. et al.: Device-Oriented Automatic Semantic Annotation in IoT. J. of Sensors 2017, 9589,064:1–9589,064:14 (2017)
20. Lubani, M. et al.: Ontology Population: Approaches and Design Aspects. J. of Information Science 45(4), 502–515 (2019)
21. Maier, M.W.: Architecting Principles for Systems-of-Systems. J. of the International Council on Systems Engineering 1(4), 267–284 (1998)
22. Matzner, M. and Scholta, H.: Process Mining Approaches to Detect Organizational Properties in CPS. In: European Conf. on Information Systems (2014)
23. Mehdiyev, N. et al.: Determination of Rule Patterns in CEP Using ML Techniques. Proc. Computer Science 61, 395–401 (2015)
24. Messager, Antoine et al.: Inferring Functional Connectivity From Time-Series of Events in Large Scale Network Deployments. Transactions on Network and Service Management 16(3), 857–870 (2019)
25. Murray, G. et al.: The Convergence of IT and OT in Critical Infrastructure. In: Proc. Australian Information Security Management Conf., pp. 149–155 (2017)
26. Ni, Z. et al.: Mining Organizational Structure from Workflow Logs. In: Proc. Int. Conf. on e-Education, Entertainment a. e-Management, pp. 222–225. IEEE (2011)
27. Noura, M. et al.: Interoperability in IoT Infrastructure: Classification, Challenges, and Future Work. In: Int. Conf. on IoT as a Service, pp. 11–18. Springer (2017)
28. Reyes-Ortiz, J. et al.: Web Services Ontology Population through Text Classification. In: Proc. Conf. on Comp. Sci. and Inf. Syst., pp. 491–495. IEEE (2016)
29. Teymourian, K. et al.: Fusion of Background Knowledge and Streams of Events. In: Proc. Int. Conf. on Distributed Event-Based Systems, pp. 302–313. ACM (2012)
30. Zhu, M. et al.: Service hyperlink: Modeling and reusing partial process knowledge by mining event dependencies among sensor data services. In: Proc. Int. Conf. on Web Services, pp. 902–905. IEEE (2017)
31. Zhuge, C. and Vaarandi, R.: Efficient Event Log Mining with LogClusterC. In: Proc. Int. Conf. on Big Data Security on Cloud, pp. 261–266. IEEE (2017)