

Enhancing Transparency of MQTT Brokers For IoT Applications Through Provenance Streams

Milan Markovic

milan.markovic@abdn.ac.uk
Computing Science, University of Aberdeen
UK

Peter Edwards

p.edwards@abdn.ac.uk
Computing Science, University of Aberdeen
UK

ABSTRACT

Systems based on Internet of Things (IoT) technologies may violate user privacy if personal data they produce and use become available to unauthorised agents. Recording provenance of IoT system behaviour may support assessment mechanisms ensuring compliance of system components with data access constraints. In this paper, we describe a prototype implementation of a provenance-enabled MQTT broker enhanced with the ability to generate provenance records describing the actual broker behaviour during message forwarding. The implementation utilises a semantic stream based approach for generating and analysing provenance data to discover message forwarding to untrusted agents. The initial evaluation demonstrates the feasibility of semantic solutions in this context.

KEYWORDS

internet of things, provenance, transparency, mqtt

1 INTRODUCTION

IoT infrastructures consist of a number of physical and virtual devices networked using a range of technologies.

The Message Queuing Telemetry Transport (MQTT) protocol [1] is a lightweight publish/subscribe messaging protocol for M2M communication. MQTT brokers are frequently used as an integral part of IoT infrastructures to provide a communication bridge between different system components [4, 9].

We argue that ensuring transparency of such brokers is a critical enabler for assessing various aspects of IoT systems such as accountability, privacy, quality and security. In this context, transparency may include information on various aspects of the deployed system such as broker configuration and the location where data processing takes place; however, it may also include information regarding data flows and agents that have access to the processed data.

For example, sensors deployed as part of an IoT system may be collecting personal data such as users' location, health information, etc. In order to ensure privacy in such settings, it is important that such data is not made available to unauthorised agents. In the MQTT context, this means that the broker should not forward messages received under a certain topic to unauthorised subscribers. To

achieve this, brokers may be extended with various authorisation mechanisms [8], which are not part of the core MQTT specification and hence not always implemented. If a broker was deployed without such a protection layer, or the layer malfunctioned, the data breach could be detected from logs generated during broker operation (e.g. information including timestamps, message payloads, clients to which the message was forwarded, etc.). However, such logs often exist as text files without any formally defined structure, which makes their processing difficult. We argue that to enhance the utility of such records, and the overall transparency of MQTT brokers, these logs should be captured in a machine-understandable format to enable intelligent audit processes to assess whether a broker is behaving as expected.

In our prototype solution presented in this paper, we apply a provenance-based approach to modelling of data describing an MQTT broker's operation during message forwarding. This produces provenance graphs that include entities (e.g. message topics), agents (e.g. subscribers receiving messages), and activities (e.g. message forwarding) as nodes and relationships between these concepts as connecting edges. In our prototype, we also implement and evaluate the linked data stream approach [2] for creating provenance streams that enable on-the-fly event detection without the need for storage of elaborate logging data which could pose scalability issues.

2 BACKGROUND

2.1 The MQTT Protocol

The MQTT protocol (version 3.1.1) is an OASIS standard¹ publish/subscribe messaging transport protocol for lightweight client-server communication in IoT infrastructures. The protocol is typically run over TCP/IP networks with small transport overhead and three basic message delivery quality assurances including the following settings: message delivered at most once (*QoS 0*); at least once (*QoS 1*); and exactly once (*QoS 2*). The protocol does not impose any restrictions on the structure of the message payload.

2.2 MQTT-Plan, EP-Plan and PROV-O

Ontologies define types of concepts and relationships between them that can be used to build semantic graphs. The graphs consist of triples where each triple has three parts: subject, predicate, and object. Each part of a triple is associated with a unique URI². Graphs can be queried by SPARQL³ queries returning instances that match query patterns (Figure 1).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

M4IoT '19, December 9–13, 2019, Davis, CA, USA

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-7028-8/19/12...\$15.00

<https://doi.org/10.1145/3366610.3368099>

¹<http://docs.oasis-open.org/mqtt/mqtt/v3.1.1/os/mqtt-v3.1.1-os.html>

²A prefix separated by : abstracts the repetitive parts of the URI. For example, ex: can represent <http://example.com#>

³<https://www.w3.org/TR/sparql11-query/>

identify provenance information that documents non-compliance (i.e. a message was forwarded to an unauthorised client) which is then saved locally as a TTL¹⁰ file. The stored provenance graphs documenting non-compliance can then be queried for relevant information, using a query of the kind presented in Listing 2.

```

PREFIX ep-plan: <https://w3id.org/ep-plan#>
PREFIX tl: <https://trustlens.org#>
PREFIX ex: <http://example.com/mqtt-plan#>
PREFIX mqtt: <http://w3id.org/mqtt-plan#>
PREFIX prov: <http://www.w3.org/ns/prov#>

SELECT Distinct ?sender ?recipient ?time
WHERE {
  ?senderVar a mqtt:Sender.
  ?senderEntity ep-plan:correspondsToVariable ?senderVar;
    prov:alternateOf ?sender.
  ?activity prov:used ?senderEntity; prov:startedAtTime ?time.
  ?result prov:wasGeneratedBy ?activity; prov:hadMember ?recipiententity.
  ?recipiententity prov:alternateOf ?recipient;
    ep-plan:correspondsToVariable ?recipientVar.
  ?recipientVar a mqtt:AffectedAgent.}
    
```

Listing 2: An example SPARQL query for inspecting provenance traces.

4 EVALUATION

In this section we describe the results of our initial evaluation of the prototype implementation. The evaluation focused on the effects of provenance mechanisms on the latency of delivered messages and time required to process semantic streams under different scenarios such as varying numbers of triples per processing window, number of messages forwarded, etc.

4.1 Set up

The evaluation was performed on a virtual server instance running Ubuntu 16.04.5 LTS with 72GB EDO DIMM and 4 CPU cores Intel(R) Xeon(R) CPU E5-2630 0 @ 2.30GHz. The same server instance was used to run the MQTT broker (started with 6GB of heap size) and the pool of MQTT clients¹¹ for publishing/receiving messages, which were implemented using the PAHO Java library¹². The stream engine used the C-SPARQL query shown in Listing 1 to monitor the provenance stream. The content of each stream window containing the provenance traces was merged with the static knowledge-base which was loaded using the C-SPARQL engine API during the broker’s startup process. The static knowledge contained the identities of trusted agents (see Listing 3).

```

<https://trustlens.org#sub-461> a <https://trustlens.org#TrustedAgent> .
<https://trustlens.org#sub-46> a <https://trustlens.org#TrustedAgent> .
    
```

Listing 3: A portion of static knowledge encoded in the TTL format. Each instance representing the ID of a known subscriber is assigned the type TrustedAgent.

Since the stream windows that the query was evaluating contained more than one provenance trace (with all elements contained within an *ep-plan:ExecutionBundle*), the query identified bundles that contained a record of the topic we were interested in observing (e.g. “/test”). Using SPARQL’s filtering function, the query also selected all triples relating to bundles which contained a record of

¹⁰<https://www.w3.org/TR/turtle/>

¹¹We experienced issues when generating more than 600 clients in a single Java application. As a result, the pool of MQTT clients was created by executing multiple jar files each creating 600 clients.

¹²<https://www.eclipse.org/paho/>

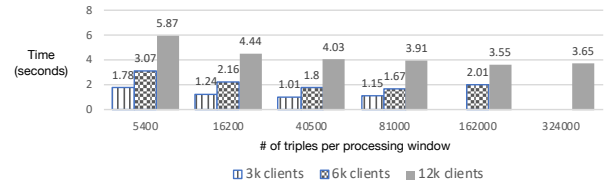


Figure 3: Time required to process a stream of provenance traces with varying processing window sizes against different volumes of subscribers.

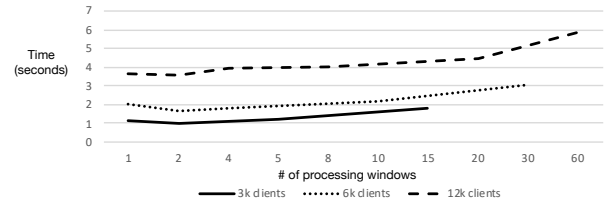


Figure 4: Time required to process a stream of provenance traces with varying numbers of processing windows against different volumes of subscribers.

an agent receiving a published message, and whose identity was not found in the static knowledge base of trusted agents.

If such triples were found, they were exported as a graph into a local file. This graph then represented a provenance trace documenting a non-compliance episode. No other provenance traces were stored by the system.

4.2 Results

Experiment 1 focused on identifying optimal sizes of processing windows (i.e. number of triples per window). We created a pool of 3k, 6k and 12k subscribers which were listening to the topic “/test” and one client publishing messages under that topic with QoS 0. We then varied the number of triples per processing window and measured the time required to process all provenance traces generated when a message was republished to all subscribers. For each setting, four messages were sent during the warm up stage. This was followed by three additional messages during which we measured the time required to process the full provenance stream. We then averaged the time from these three observations (the same approach was applied to other experiments described in this section). Figure 3 and 4 show that the system performed best when the processing load was split in 1-2 processing windows¹³. As illustrated by Figure 4 there seems to be a small saving if the overall processing load is split into two windows.

Experiment 2 evaluated whether our implementation of provenance mechanisms had any significant impact on the baseline performance of the Moquette broker. We created different republishing loads on the broker and measured the latency at which the messages were received by clients. In this experiment we again created

¹³Note that with 3k subscribers and window size set to 81k triples the system only processed one window (i.e. the size of a single provenance trace describing message forwarding to one subscriber is 27 triples)

a network of a single publisher publishing under a single topic to a pool of subscribers. Each message was associated with a timestamp, which was generated when the message was sent, and each client then calculated the latency based on the time when the message was received. Table 1 shows a very small increase of 40 - 50 ms in average message latency for 3k and 6k subscribers.

		Number of subscribers	1200	3000	6000
Without provenance	MIN		29	28	36
	MAX		1039	1064	1105
	AVG		532	545	552
With provenance	MIN		35	31	35
	MAX		1045	1106	1119
	AVG		531	588	603

Table 1: Observed message latency in ms.

In *Experiment 1*, for every test we loaded a static knowledge base containing 12k of triples and each test generated one provenance trace of non-compliance identifying a single untrusted agent. In *Experiment 3*, we were interested in evaluating whether the size of the static knowledge base, and varying populations of trusted and untrusted agents present in the stream had any impact on processing time. We created a pool of 6k subscribers and set the window size to 81k triples (based on the results from *Experiment 1*). We then varied the size of the static knowledge containing information about trusted agents and the populations of untrusted agents receiving messages.

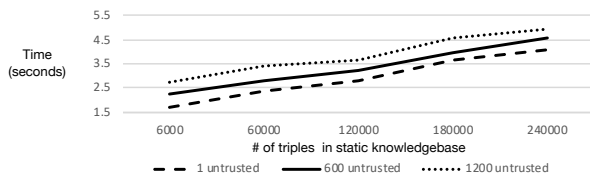


Figure 5: Time required to process a stream of provenance traces with varying numbers of untrusted agents and varying static knowledgebase size.

Figure 5 shows the results, indicating that increasing the size of the static knowledge base indeed has a negative impact on the time required to process the provenance stream. However, the function of the growth of processing time appears to be free of any exponential increases that could hinder the scalability of the solution. Figure 5 also shows that the processing time is impacted by the number of untrusted agents that are detected. This could be explained by the higher demand for resources from file writing processes that have to write larger files describing saved provenance traces. In future, this could perhaps be improved by writing into an in-memory triple store rather than onto a disk.

4.3 Limitations

Currently, the proposed approach has a number of limitations. For example, it assumes that information captured in the static knowledge base (i.e. the list of trusted agents) can be linked to the ID of the MQTT client. An alternative approach could be to calculate the trustworthiness of a client on the fly based on some characteristics

included in the provenance trace that can be observed by the broker (e.g. sender's location).

Another limitation is related to the case where the size of the stream processing window is larger than the sum of the triples contained by the generated provenance traces. Such a window would not be evaluated until additional provenance traces had filled the required window size. This could be mitigated by altering the C-SPARQL engine implementation to force such windows to be evaluated after a certain time period (e.g. 500ms).

Finally, using the same server instance to run the broker and the client pool would mean more CPU demand which could have an effect on the observed performance results.

5 CONCLUSIONS & FUTURE WORK

In this paper, we presented a partial implementation of a semantic provenance-based transparency mechanism within a third party MQTT broker. Our initial evaluation results suggest that implementing semantic solutions within middleware IoT components offers a feasible and potentially scalable approach to enhance their auditability.

In our future work, we will focus on extending the provenance mechanisms to cover other parts of a broker's behaviour (e.g. handling connection requests) as well as updating the MQTT-Plan ontology to align with the recently introduced version 5 of the MQTT protocol¹⁴. We will also explore the opportunities for recording additional plan metadata provided by the EP-Plan ontology and integration of MQTT provenance in the wider IoT system context. We also plan to test the separation of the stream processing engine from the broker in a distributed network set up.

ACKNOWLEDGMENTS

The work described here was funded by the award made by the RCUK Digital Economy programme to the University of Aberdeen (EP/N028074/1).

REFERENCES

- [1] BANKS, A., AND GUPTA, R. Mqtt version 3.1.1. Tech. rep., OASIS Standard. <http://docs.oasis-open.org/mqtt/mqtt/v3.1.1/os/mqtt-v3.1.1-os.html>, 29 October 2014. Latest version: <http://docs.oasis-open.org/mqtt/mqtt/v3.1.1/mqtt-v3.1.1.html>.
- [2] BARBIERI, D. F., BRAGA, D., CERI, S., VALLE, D. E., AND GROSSNIKLAUS, M. Querying rdf streams with c-sparql. *SIGMOD Rec.* 39, 1 (2010).
- [3] GARJO, D., AND GIL, Y. Augmenting prov with plans in p-plan: scientific processes as linked data. CEUR Workshop Proceedings.
- [4] LAMPKIN, V., LEONG, W. T., OLIVERA, L., RAWAT, S., SUBRAHMANYAM, N., XIANG, R., KALLAS, G., KRISHNA, N., FASSMANN, S., KEEN, M., ET AL. *Building smarter planet solutions with mqtt and ibm websphere mq telemetry*. IBM Redbooks, 2012.
- [5] LEBO, T., SAHOO, S., AND MCGUINNESS, D. PROV-O: The PROV ontology. Tech. rep., April 2013.
- [6] MARKOVIC, M., CORSAR, D., ASIF, W., EDWARDS, P., AND RAJARAJAN, M. Towards transparency of iot message brokers. In *International Provenance and Annotation Workshop (2018)*, Springer, pp. 200–203.
- [7] MARKOVIC, M., GARJO, D., EDWARDS, P., AND VASCONCELOS, W. Semantic modelling of plans and execution traces for enhancing transparency of iot systems. In *2019 Sixth International Conference on Internet of Things: Systems, Management and Security (IOTSMS) (2019)*, IEEE.
- [8] NIRUNTASUKRAT, A., ISSARIYAPAT, C., PONGPAIBOOL, P., MEESUBLAK, K., AIUM-SUPUGGUL, P., AND PANYA, A. Authorization mechanism for mqtt-based internet of things. In *2016 IEEE International Conference on Communications Workshops (ICC) (2016)*, IEEE, pp. 290–295.
- [9] SONI, D., AND MAKWANA, A. A survey on mqtt: a protocol of internet of things (iot). In *International Conference On Telecommunication, Power Analysis And Computing Techniques (ICTPACT-2017) (2017)*.

¹⁴<https://docs.oasis-open.org/mqtt/mqtt/v5.0/mqtt-v5.0.html>