# LUCON: Data Flow Control for Message-Based IoT Systems

Julian Schütte, Gerd Stefan Brost

Fraunhofer AISEC, Germany

{julian.schuette,gerd.brost}@aisec.fraunhofer.de

*Abstract*—Today's emerging Industrial Internet of Things (IIoT) scenarios are characterized by the exchange of data between services across enterprises. Traditional access and usage control mechanisms are only able to determine *if* data may be used by a subject, but lack an understanding of *how* it may be used. The ability to control the way how data is processed is however crucial for enterprises to guarantee (and provide evidence of) compliant processing of critical data, as well as for users who need to control if their private data may be analyzed or linked with additional information – a major concern in IoT applications processing personal information. In this paper, we introduce LUCON, a data-centric security policy framework for distributed systems that considers data flows by controlling how messages may be routed across services and how they are combined and processed. LUCON policies prevent information leaks, bind data usage to obligations, and enforce data flows across services. Policy enforcement is based on a dynamic taint analysis at runtime and an upfront static verification of message routes against policies. We discuss the semantics of these two complementing enforcement models and illustrate how LUCON policies are compiled from a simple policy language into a first-order logic representation. We demonstrate the practical application of LUCON in a real-world IoT middleware and discuss its integration into Apache Camel. Finally, we evaluate the runtime impact of LUCON and discuss performance and scalability aspects.

## I. Introduction

While IoT systems in general create undeniable benefits in areas like health care, home automation, manufacturing, logistics, and mobility, it is also obvious that existing threats to the integrity of business processes and the privacy of users intensify with an increasing degree of distribution, amount of endpoints and trust domains. A paramount challenge for data owners is to control the way how their data is processed and combined with data from other sources, and how it is published to untrusted third parties.

Today, Industrial IoT systems are characterized by data flowing from sensors to services and applications and possibly back to actuating devices. These data flows span several physical platforms, including resource-constrained sensors, mobile devices, and cloud backends. In contrast to traditional enterprise systems, modern distributed IoT systems typically span several "trust domains", i.e. within a single application, data is processed by services under different authoritative controls.

In addition, privacy and security are interleaved since sensor data may contain personal data of employees (e.g.,

a machine operator) or private users (e.g., the owner of a home automation solution).

Traditional access control cannot cope with this challenge – it merely aims at controlling actions of subjects to resources (e.g., a "read" request from a user to a file). In that sense, traditional access control is *resource-centric* and unaware of the actual processing of data, as access to resources is only controlled at a specific point in time without further control on how it is used. Further, typical access control languages like XACML provide means to describe resources, but do not allow to write rules referring to classes of data that provided by these resources. It is impossible to state that only certain information may be retrieved from an endpoint, while some other information must not be published by the same endpoint.

An extension of access control is *usage control* which has been introduced in the early 2000s [16] and has been subject to intensive research in the following decade [10], [12], [13]. Usage control extends access control by the dimension of time and is able to continuously monitor and control the usage of resources such as files or services by subjects. However, it is mostly still resource-centric as it only decides access requests to resources in the course of time, but does not control how sensitive information is processed and combined. ABAC languages like XACML 3.0 are moving in that direction by supporting the notion of *obligations* which must be fulfilled by the subject. The outcome of an obligation does however not influence the policy decision anymore and is out of the semantics of XACML.

In modern data-centric systems, the concept of resource-centric protection does not apply anymore. As a consequence, traditional usage control models fall short of enforcing requirements of data owners. With a growing number of data sources in the form of sensors from different owners and cloud-based data analytics services, the challenge is not anymore to control the usage of a single resource, but to express constraints on how *data objects* (messages) may be processed . Figure 1 illustrates a typical scenario: sensors in a production facility measure parameters of the production process like flow rate and temperature of various liquids. These measurements are essential for controlling the production process, but they are also interesting for analytics applications from third parties. Knowing these measurements helps manufacturers
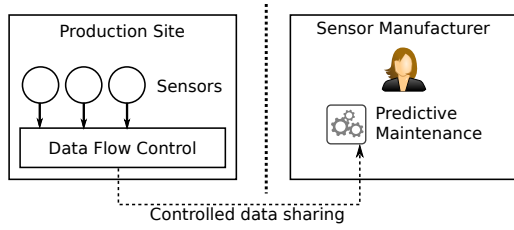
Figure 1. Data flow control for predictive maintenance

developing "predictive maintenance services" such as the detection of sensor drifts or an approaching end of life of their hardware.. However, up to date, these scenarios are hardly possible due to the sensitive nature of raw sensor data, from which trade secrets like recipes, production processes and capacities can immediately be derived.

So, controlling access to resources (in this case, sensors) or their data is not sufficient. Rather, it is necessary to control the flow of data, i.e. the way how it is combined, processed and shared with different endpoints. In the example from Figure 1, it must be guaranteed that sensor manufacturers only get data from specific sensors. This data must pre-processed in a way that allows them to run their analytics, but not to reverse-engineer the production process or any other (including privacy related) data.

In this paper we introduce LUCON, a policy framework for data flow control (DFC) in distributed systems, which provides a runtime monitor for dynamic enforcement of DFC policies and an upfront static verification of message routes against policies. As discussed in [25], we model all data exchanges as message flows. Data flow policies are based on a formal system model and an operational semantics of the enforcement in message routes, which pre-configure possible data flows. After introducing the formal foundation of LUCON we show how it is implemented in a real-world messaging system and supports both static and dynamic enforcement. The benefit of an upfront static enforcement is that users can analyze potential, possibly counter-intuitive violations of their policy, while in a dynamic enforcement only concrete violations of policies will be prevented. In our policy framework, the result of a policy enforcement is either a simple cancellation of a message flow or the execution of an obligation, i.e. an action which must be taken by the enforcement component to fulfill the policy. By means of a prototype implementation of the LUCON framework and its integration into the Apache Camel messaging router, we show that the performance of LUCON is suited for large-scale productive applications.

## II. RELATED WORK

Usage control has been subject to extensive research for quite a while [24]. While several models have been proposed, the most prominent one is $UCON_{ABC}$, originally introduced by Park and Sandhu [16]. It comprises Authorizations (A), oBligations (B), and Conditions (C), referring to attributes of subjects and resources. Attributes

are mutable (e.g., they can change over time) and the continuity of access decisions is formalized. In this way, UCON A, B and C can be defined to be evaluated before (pre) or during usage (on). The model has undergone different extensions in the course of time. As an example, [12] incorporated post-obligations. $UCON_{ABC}$ does not dictate how to design a specific architecture and mechanisms for usage control, but stays abstract in that manner. Other approaches focus on specific languages rather than abstract models, such as the Obligation Specification Language (OSL) [10].

Much has been done in the area of formalization of usage control policies and the formal analysis of their properties. In [27], a formalization of $UCON_{ABC}$ in Lamport's Temporal Logic of Actions (TLA) is given, in [2], [1] Basin et al. give an approach on analyzing usage control policies formalized in first-order temporal logic (MFOTL). In [22], a Linear Time Logic (LTL) dialect is used for the sake of analyzing policies, and in [8] an analysis of dynamically changing usage control policies is described, based on Action Computation Tree Logic (ACTL). Our work is based on this research, but the approach is more specific and focuses on the application of usage control to data processing only.

The concept of enforcing data flow control in decentralized systems has already been introduced by Myers et al. [14]. Their understanding of controlling information flows refers to preserving secrecy and integrity properties of classified documents – an approach that pursues the enforcement of traditional information classification systems such as the Bell LaPadula model (*no read up, no write down*) for secrecy and the Biba model [4] (*no read down, no write up*) for integrity. Myers proposes a label-based approach to mark data sets and to prevent information leakage by annotating existing programming languages. We generalize this concept by concentrating on information flow between components (services), that have no built in mechanisms for supporting external labels. The enforcement of usage control policies is a central challenge, as it requires system-specific implementations and trust relationships related components. Trustworthy system architectures for usage control enforcement have been proposed in [28], which allow usage control policies at the level of system calls, given that the trustworthiness of the enforcement point can be attested using hardware-based mechanisms. Our approach does not focus at usage control enforcement on remote platforms, but rather on a specific enforcement mechanism for data processing. Nevertheless, it could be combined with techniques from [21] or [28] in case the enforcement would have to take place on remote hosts.

Closer related to our work is [9], which introduces the idea of using data flow tracing at the level of system calls in order to enforce usage control policies. The authors show that based on an underlying data flow model, more realistic

and expressive policy rules can be written, referring to states of a data flow system, rather than specific sequences of events. In [20], this approach is extended by tracing messages in the X11 environment, specifically copy & paste actions on sensitive data which is either blocked or replaced by meaningless data in case a policy is violated. Similar to [9], [22], we understand usage control as enforcing conditions in data flows.

Fine granular data flow tracking for databases has been done by applying taint tracking in [6] for specific applications. A similar approach was followed by [5], providing an API that allows to introduce taint tracking for legacy web applications without major code changes. This is also based on taint tracking at database level and uses hooks that are placed in legacy code to enable security policy enforcement.

Pasquier et al. proposed CamFlow [18], an end-to-end information flow control enforcement system for cloud systems based on the implementation of Linux Security Modules as enforcement points. Thus, CamFlow is tightly integrated into the operating system via a custom Linux Security Module and considers information flows between processes in an OS. This work has been continued in [17] with a focus on DETA (Declassify, Endorse, Transform, Authorize) policies. Apart from the fact that CamFlow is an operating-system-level mechanism while our approach mainly addresses message buses in a distributed system, CamFlow proposes fixed security rules for secrecy and integrity, following traditional information classification concepts. Our system, in contrast, proposes a more generic labeling mechanism that allows to create any information class, track its usage in the system and write policies on how to handle read and write accesses to it.

## III. OVERVIEW OF OUR APPROACH

LUCON is a policy framework to enforce secure data flows in message-based systems – typically in IoT architectures. Its design goals are a reasonable low runtime overhead, a formal semantics and support of authors in writing flawless policies for existing message routes. It is comprised of the following components:

- the definition of a policy language, its implementation in Eclipse XText, and its compilation into a first-order logic representation
- a runtime evaluation of policies based on the first-order logic representation of policies, following a message tainting approach with a formal execution semantics of policy-controlled message routes
- a static model checking of message routes against policies in the first-order logic representation, including a compilation of Apache Camel message routes into Prolog programs

### A. Policy language

The motivation for a policy language is to separate the specification of security requirements on data flows

from the actual messaging system. Existing DFC systems like [18], [15], [19] enforce predefined flows between security classes. However, in practice users have diverse and application-specific requirements which cannot be hard coded into a generic distributed middleware. Early research on information flow control has proposed various information classification models such as Chinese Wall, Biba or Bell LaPadula, where each serves one specific requirement (such as either integrity or secrecy for lattice-based information classifications), but they are not necessarily compatible with one another. Thus, instead of hard coding valid flows into the system, we rather aim for a simple domain specific language (DSL) to allow the user to define custom policies. The DSL compiles into a logic representation in Prolog – a programming language based on Horn clauses which is a decidable subset of first-order logic. The benefit of that is that compiled policies have a formal foundation and can be used to model-check requirements against message routes, but at the same time can be efficiently evaluated at runtime so that performance impact on a productive system remains low.

### B. Runtime enforcement

Runtime enforcement implements a dynamic taint-style analysis and thus leans towards the permissive end of possible data flow enforcement strategies, as we will discuss in section V. The aim of LUCON is to prevent messages from violating the policy, but not to prevent any information leaks over side channels. Some data flow systems proposed in the past [23] apply a stricter strategy and block even information leaks over side-channels, for example in cases in which the attacker can learn private information by observing the control flow or exceptional terminations of a message route. However, these approaches assume that the attacker knows the exact control flow specification (i.e., the message routes), which is not the case in our model. Second, implicit information leaks are not intuitive for the user and sudden cancellation of a route at runtime may not be expected behavior. In our taint-style approach we mark messages with an initial set of labels as soon as they are created. Labels are transported along with messages and possibly altered when the message is processed by a service. The mechanism for the modification of message labels is called the *taint propagation logic* and determines how the security and privacy properties of a message change as it is processed and merged by services. In our approach, the definition of the taint propagation logic is part of the policy. This allows our runtime monitor to query the compiled policy for the changes which should be made to message labels, as well as for the actual data flow policy.

### C. Static model-checking

Runtime enforcement is tuned to be fast and will prevent messages from leaking information by blocking or modifying them just before the data leak would occur.

For users, it is however important to analyze if and under which circumstances their message system would run into potential data leaks so they can verify that message routes will not be unexpectedly terminated by the policy framework. Further, LUCON will provide evidence that message routes fulfill the security requirements, which is important information for audit and compliance purposes. This is achieved by translating message routes into a first-order logic representation, analog to the representation of compiled policies. The logic model allows to verify route definitions against policies so that users can check upfront whether routes are applicable at all under a certain policy, whether only specific executions paths may violate the policy, or whether a route is fully compliant with a policy. In case of potential policy violations, LUCON will generate an example of a message flow violating the policy.

## IV. SYSTEM MODEL

We first set the common ground for the abstract type of system that is addressed by our policy framework. In practice, LUCON runs in any message-based IoT system, but for the remainder of this paper we establish an understanding of the terms and concepts that are relevant in those systems.

The system is based on *services* which communicate via *messages*. A service accepts a set of input messages, operates on their content and emits a set of output messages. Each service is under control of a trust domain and as messages are sent from one service to another, they may cross domain boundaries. The ability of a user to define and apply policies is limited to their own trust domain, i.e. policies of a user can only control messages within their domain. With respect to enforcement of policies, however, it is still possible that a user retrieves an assertion of a successful enforcement from a remote domain – either by establishing trust at a technical level (e.g., by remote attestation) or by retrieving evidence of the enforcement (e.g. by observing expected side-effects).

Both, We denote these sets of predicates as message labels $\mathcal{L}$ and service *properties* $\mathcal{P}$, respectively.

*a) Message Labels:* Message labels $\mathcal{L}$ classify a message in terms of its data source or secrecy level and may be partially ordered. For instance, a message $m$ can be labeled as $\mathcal{L}_m = \{\texttt{classification(top\_secret)}\}$ (1-ary predicate) or $\mathcal{L}_m = \{\texttt{personal\_data}\}$ (0-ary predicate). The specific predicates are not determined by the model but rather by its instantiation in a specific application.

*b) Service Properties:* Service properties $\mathcal{P}$ are used to describe services. For example, a service which stores data in a database can be assigned the predicates $\mathcal{P} = \{\texttt{persist}\}$ (0-ary predicate) or $\mathcal{P} = \{\texttt{persist(jdbc://localhost/...)}\}$ (1-ary predicate).

*c) Message Routes:* The interaction between services is defined as an Enterprise Integration Pattern (EIP) [11] in form of a message *route*. We consider a route as a
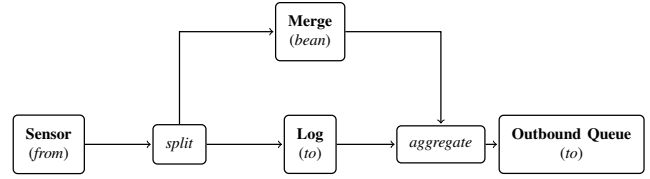


Figure 2. A message route over several services

non-while-looping program, i.e. a sequence of numbered statements which either call external services, assign values to variables, or control execution of the next statement. Note that excluding while loops from our route definition is a limitation compared to the expressiveness of real-world turing-complete message routers like Apache Camel or Spring Integrations, which do in fact allow the construction of while loops in EIPs like *Dynamic Router*[1]. The reason we chose to exclude while-loops is that it turns the static route verification into a decidable problem, while in practice while-loops are rarely used in EIPs and e.g. discouraged by Apache Camel[2]. Routes support variables in two scopes: global and message-scoped. Global variables are available across all executions of a route, while message-scoped variables get appended to the message object and are transported along with it. Branching statements refer to conditions over variables and fork the control flow into several branches, just like conditions in a program.

Accordingly, the set of supported statements comprises variable assignments (set-*-prop), control flow modification (choice), message manipulation (split, aggregate, bean), and service invocation (from, to). The simplified grammar is given in the following listing and Figure 2 depicts an example route.

$$
\begin{array}{rcl}
stmt & := & assign\text{-}msg \mid assign\text{-}env \mid from \\
      &   & \mid to \mid choice \mid split \mid aggregate \\
assign\text{-}msg & := & \texttt{set-msg-prop}\ \text{var} := expr \\
assign\text{-}env & := & \texttt{set-env-prop}\ \text{var} := expr \\
from & := & \texttt{from}(service) \\
to & := & \texttt{to}(service) \\
choice & := & \texttt{when}\ expr\ \texttt{then goto}\ v \\
       &   & \texttt{otherwise goto}\ v \\
split & := & \texttt{split}\ expr \\
aggregate & := & \texttt{aggregate}\ expr \\
expr & := & \text{n-ary Prolog predicate} \\
v & := & \text{Statement number} \\
service & := & \text{Service name}
\end{array}
$$

To model execution of a route, we further introduce some *execution contexts* which represent the current state of the execution: $\Sigma$ maps statement numbers to statements. $\mu_m$ holds the message-scoped variables and maps each variable of a message $m$ to its value. A global map $\Delta$ assigns global variable names to their values. Further, a program counter pc holds the number of the currently executed

[1]http://www.enterpriseintegrationpatterns.com/patterns/messaging/DynamicRouter.html

[2]cf. http://camel.apache.org/loop.html

| $\tau$ | Maps a message to its taint state, e.g. $\tau[m \leftarrow 1] \Rightarrow 1 = \tau[m]$ |
|---|---|
| $\Sigma$ | Maps a statement number to a statement |
| $\mu_m$ | Maps variables of message $m$ to their value |
| $\Delta$ | Maps global variable names to their current value |
| pc | Number of the currently executed statement |
| $\iota$ | Number of the next statement |

statement and an instruction pointer $\iota$ the number of the next statement.

Table I summarizes these execution contexts.

*d) Trust Domains:* Services and routes reside in *trust domains*. A trust domain is controlled by a single authority that can create and update route definitions and it is assumed that services within a trust domain behave correctly in terms of propagating message labels. It is important to note that we assume that route definitions are not known outside of the trust domain. If this assumption would not hold, information from messages may leak if an attacker is able to observe the control flow, i.e. the execution of a route.

## V. HYBRID INFORMATION FLOW CONTROL

LUCON takes a hybrid approach on information flow control by combining a dynamic and a static component: The dynamic policy enforcement is tuned for efficiency and to limit delays by policy evaluations at runtime. It is based on a *taint-style analysis* that prevents explicit information leaks under the assumption that implicit leaks (for example by control flow observation or untrusted misbehaving services) do not occur. A static, upfront model-checking verifies message routes against policies and informs the user about potential policy violations. Here, runtime performance does not play a role, but rather completeness of the verification and the generation of understandable counterexamples so as to either guarantee that message routes are free of policy violations (e.g., for audit purposes) or to support policy authors in fixing potential flaws.

### A. Static vs. Dynamic Data Flow Control

Data flow research dates back to the seventies when Denning [7] proposed a lattice-based organization of security classifications to mathematically formulate constraints on information flows – a formal foundation for the Bell-LaPadula secrecy [3] and Biba [4] integrity models. Since then, various data flow control systems have been proposed, either relying on static checking of system configurations against information flow rules, or on dynamic enforcement of data flow constraints at runtime [19], [18], [15], [26]. These systems typically enforce secrecy by preventing information leaks over explicit flows and partly over implicit flows. Explicit flows refer to leakage of information directly into publicly readable sinks, whereas implicit leaks refer to side-channel leaks via control flow or termination. Sabelfeld et al. show in [23] that dynamic enforcement and the classic Denning-style static flow control can only

achieve *termination-insensitive non-interference*, i.e. they can prevent information leaks via observation of the control flow by an attacker, but not via observation of route termination. Taint analysis, as we adopt for our dynamic runtime enforcement, provides even weaker guarantees as it allows control-flow leaks in some cases. To illustrate this, let us consider the following route, which we denote as a simple program, for the sake of readability.

```
1 tainted := ...;  // Taint label set
2 public := 1;
3 tmp := 0;
4 if tainted then
5    tmp := 1;
6 if tmp != 1 then
7    public := 0
```

In this example, sensitive information is written into a variable `tainted`, which will consequently be marked with a taint label (line 1). As the value of `tainted` is never assigned to any other variable, the taint flag is not propagated. Variable `public` is written into a public data sink and thus leaks its content. As can be seen from the example, the value of variable `public` is equal to `tainted` in all possible execution paths, although it is never explicitly assigned. Consequently the route leaks tainted information to a public data sink and a classic taint analysis is not able to detect this leak.

A Denning-style type system would behave differently and manage a stack of global security contexts. In line 4, when a "secure" condition is evaluated, a new *secure* item would be added to the stack. When execution enters line 5, the system would notice that a non-secure variable is written within a secure context and would terminate execution immediately, thereby preventing the information leak. Denning-style type systems are thus more strict and prevent information leaks even under the assumption that the attacker know the control flow specification (the message route) and is able to observe control flow at runtime. However, in the type of distributed system we address, Denning-style flow control is less appropriate, as it introduces a variety of practical issues: first, the assumptions are overly strict. Attackers might be able to observe parts of the control flow, e.g. by hosting a service which is used by a message route, but they do not know the control flow specification (the message routes) nor can they learn it by globally observing the control flows of domain. Second, in real-life message routes, all operations would have to be considered as write operations, as they are typically realized by some implementation which is not further known to the data flow control engine. If all operations are writes, however, a single access to a non-tainted variable will lead to termination of the route. This is unexpected for the user, at best, and will render the system dysfunctional.

As a consequence, LUCON adopts a dynamic taint approach which is efficient at runtime and prevents explicit leaks of information, i.e. it prevents routes from processing data in an undesired way, assuming that an attacker cannot retrieve the route specification and globally observe

the control flow. Complementary to runtime enforcement, LUCON provides an upfront static model checking to verify message routes against policies. In the approach we describe herein, the model checking follows the same taint-style semantics as the dynamic enforcement, but in general the models do not have to be equal. For instance, it would be possible to statically verify routes in a stricter control-flow- and termination-sensitive model to identify even theoretical information leaks, while still running dynamic enforcement in the more realistic and relaxed taint-style model.

### B. Dynamic Taint-Style Enforcement

The basic idea of the taint-style dynamic flow enforcement is to assign a set of taint labels to messages when they enter the system and to modify the taint labels as messages are processed by services. Whenever a message is about to be sent to an external service, the policy is consulted to check whether a respectively marked message may enter this specific service.

Different from other information and data flow control systems, LUCON does not dictate a set or lattice of taint labels, but rather allows to assign any set of labels to messages. Taint labels are represented as first-order logic predicates and any rule over these predicates can be declared to construct lattices, hierarchies or any other inference of labels. Assignment of labels to messages is controlled by the taint propagation logic, which is part of the policy. In fact, every service description in the policy may include two *label transformation functions* $\mathcal{L}^-(\cdot)$ and $\mathcal{L}^+(\cdot)$ that determine which labels will be removed and added to a message, respectively. To denote the semantics of a route with taint tracking enabled, we introduce an additional context $\tau$ that maps variables to the set of taint labels assigned to them. That is, $\tau_\Delta$ denotes the taint states of global variables and $\tau_m$ denotes the taint labels assigned to a message $m$. Individual variables of a message cannot be tainted, rather the whole message will be marked.

The operational semantics of a taint-controlled route execution is given in the appendix in Figure 7. Inference rules are written as

$$\frac{Computation}{\langle Current\ state\rangle, Stmt \to \langle Next\ state\rangle}$$

where $Current\ state$ and $Next\ state$ are written as tuples $\langle \tau, \Sigma, \mu_m, \Delta, \mathrm{pc}, \iota\rangle$ and denote the system state before and after execution of the statement $stmt$, respectively. $Computation$ denotes the actual computation on the system state which is applied by executing $Stmt$. The notation of computations makes use of expressions in the form $\mu_m, \Delta \vdash e \Downarrow v$, which means that an expression $e$ evaluates to value $v$ in the context denoted by messages properties $\mu_m$ and system variables $\Delta$.

Statements refer to operations of typical enterprise integration patterns (EIP), as used by Apache Camel[3]. The formal semantics covers far from all Camel EIPs but is focused on the statements relevant for information flows.

The FROM statement reads data from a service endpoint and TO and BEAN forward it to an external service or an internal processing bean, respectively (a component that may modify the message). CHOICE denotes a branch in the control flow and is similar to an if-then-else-statement in a normal program. With SPLIT, a message can be split by an expression into multiple messages which are processed in parallel and can be joined again by AGGREGATE. Statements ASSIGN-MSG-PROP and ASSIGN-ENV-PROP set message-scoped and global-scoped variables to the value of a given expression. Variables are only visible within the message routing engine and not delivered to actual services, therefore they do not affect the taint state of a message. Rather, the only statements affecting the taint state $\tau$ are FROM, TO, BEAN, SPLIT, and AGGREGATE. When a message is created by FROM, it is assigned the taint labels determined by the taint policy $\mathcal{L}^+$. When that message is forwarded to any other service, the taint labels according to $\mathcal{L}^-$ are removed and the ones determined by $\mathcal{L}^+$ are added. When a message is split, all resulting messages have the same taint labels as the original one and when messages are merged, the resulting message is tainted with the union of all individual taint labels.

### C. Static Model-Checking of Data Flows

Dynamic taint tracking at runtime is sound under the assumption that the attacker does not know the message route definition, i.e. the control flow, but it is not complete, in the sense that it will only detect actual data leaks as they occur and not guarantee that a message route is free of data leaks in general. However, as it is important for users to know if a route may be interrupted by a policy, we use static model-checking to verify routes against policies.

For this purpose, we compile routes into Prolog, i.e. the same logic representation as policies. In Prolog, a route is represented as a directed acyclic graph, where each node represents one statement and edges refer to transitions between statements. Predicate `stmt` defines a statement and `succ(A,B)` defines `B` as a successor of `A`. This way the example route from Figure 2 can be written as the following (simplified) Prolog program:
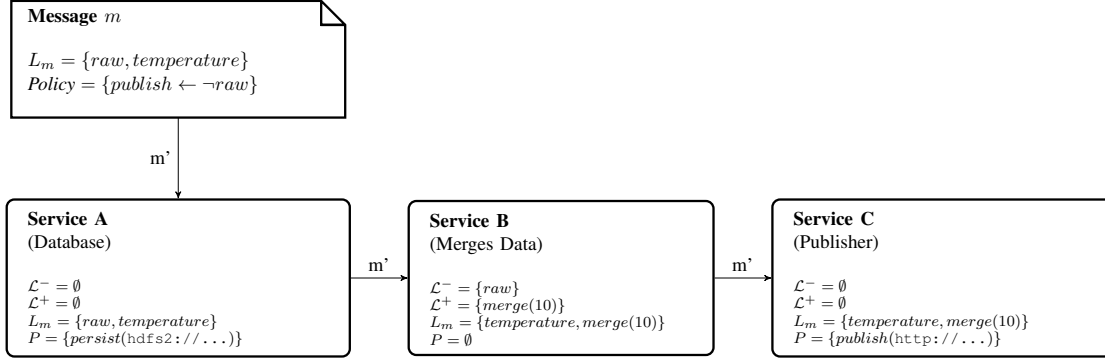
```
1 stmt(sensor).          7 succ(sensor,split).
2 stmt(split).           8 succ(split,log).
3 stmt(log).             9 succ(split,merge).
4 stmt(merge).          10 succ(merge,aggr).
5 stmt(aggr).           11 succ(log,aggr).
6 stmt(mqueue).         12 succ(aggr,mqueue).
```

Policies are likewise compiled into Prolog and determine valid and invalid flows in terms of allowed and forbidden labels entering services. Message routes can then be checked against policies by respectively exploring all paths in the

---

Figure 3. Message $m$ with taint labels $L_m$ sent along services with properties $P$

**Message** $m$

$L_m = \{raw, temperature\}$
$Policy = \{publish \leftarrow \neg raw\}$

m'

**Service A**
(Database)

$\mathcal{L}^- = \emptyset$
$\mathcal{L}^+ = \emptyset$
$L_m = \{raw, temperature\}$
$P = \{persist(\texttt{hdfs2://}...)\}$

m'

**Service B**
(Merges Data)

$\mathcal{L}^- = \{raw\}$
$\mathcal{L}^+ = \{merge(10)\}$
$L_m = \{temperature, merge(10)\}$
$P = \emptyset$

m'

**Service C**
(Publisher)

$\mathcal{L}^- = \emptyset$
$\mathcal{L}^+ = \emptyset$
$L_m = \{temperature, merge(10)\}$
$P = \{publish(\texttt{http://}...)\}$

graph which violate the policy. Each solution to the query is one counterexample of a possible data flow in a message route which does not comply with the policy. Listing 1 shows the output displayed when a route violates a data flow policy.

Listing 1. Proof of a message route violating a policy

```
1  Route Sensor_Messaging is invalid because
2  service Outbound_Queue may receive label(s) [raw].
3  This is forbidden by rule dontPublishRaw
4
5  Example flows violating policy follow:
6  |-- sensor creates message labeled [raw]
7  |-- split receives message labeled [raw]
8  |-- log receives message labeled [raw]
9  |-- aggr receives message labeled [raw]
10 |-- mqueue receives message labeled [raw]
11 |-- fail
```

## VI. THE LUCON POLICY LANGUAGE

So far, we described how LUCON controls data flows in terms of abstract models, which provide the formal foundation of our policies. To be of any practical use, the framework must allow to write policies in a language that is easy to understand and supports the user in writing correct policies.

The LUCON policy language is a domain specific language (DSL) which serves two purposes: first, it comprises the actual data flow control rules, determining valid and invalid data flows and possibly binding them to obligations. Second, it defines labels and describes services in terms of properties, capabilities and their taint propagation logic $\mathcal{L}^-$ and $\mathcal{L}^+$.

We define a grammar for the LUCON DSL in Eclipse XText. XText is a language creation framework that automatically creates lexer and parser from a context-free LL(*) grammar, along with IDE editors with syntax highlighting, auto-completion, and error checking.

The following is a simplified version of the LUCON DSL grammar. It defines the main concepts *service* and *rule*, which represent the service-specific taint propagation and the actual data flow rules. The effect of a rule is represented by a *decision* that determines whether a message may be

```
flow_rule {
    id dontPublishRaw
    when {
        endpoint "http[s]?://.+"
    }
    receives raw
    decide  drop
            require log("Preventing data leak. ", message)
    otherwise error
}
```

Figure 4. LUCON DSL rule in Eclipse IDE

passed on or must be dropped. Optionally, a decision can be bound to an *obligation*. Obligations are actions which must be executed successfully before the actual decision is enforced. If the execution of an obligation fails or the respective obligation is not supported by the system, the alternative decision stated by `otherwise` is taken.

$$
\begin{aligned}
policy & := rule* \mid service* \\
service & := \texttt{service \{} \\
& \qquad \texttt{id } atom \\
& \qquad \texttt{endpoint } url \\
& \qquad \texttt{(properties } term\texttt{+)?} \\
& \qquad \texttt{(capabilities } term\texttt{+)? \}} \\
rule & := \texttt{flow\_rule \{} \\
& \qquad \texttt{when } s \texttt{ receives } atom \\
& \qquad \texttt{decide } decision \texttt{ \}} \\
term & := \text{Prolog term} \\
atom & := \text{Prolog atom} \\
url & := \text{Endpoint URL of a service} \\
decision & := effect \; (obligation)* \\
effect & := \texttt{allow} \mid \texttt{drop} \mid \texttt{error} \\
obligation & := \texttt{require } term \; (\texttt{otherwise } term)? \\
s & := \text{Reference to a } service
\end{aligned}
$$

Figure 4 shows a policy from the example scenario in section I. It includes an inline definition of the services it refers to – in this example simply all services with an http(s) endpoint. The message label `raw` is stated as an atom (i.e., a 0-ary predicate) and marks raw sensor data. If that label has not been removed along the message route by some service that merges or blinds raw data records, the rule is triggered and will drop the message before it

enters the respective endpoint. Before, the event will be logged, whereas `log` refers to a Java function which can be called from within the policy decision point and `message` refers to a predefined variable holding the content of the message. In case the execution of the obligation fails, the rule's effect is `error` which exceptionally terminates the message route.

LUCON policies are compiled into Prolog programs using the Xtend code generation framework, so that users only deal with the high level DSL, while the enforcement engine operates on the formal representation in Prolog. The representation of a policy in a logic model further allows reasoning over the policy itself to detect conflicting or incomplete rules and provides the basis for the aforementioned static model checking of message routes against data flow policies. Listing 2 shows the Prolog representation of the rule from Figure 4.

Listing 2. Prolog representation of a rule

```
1 regex(A,B,C) :- class("j.u.r.Pattern")
2                     <- matches(A,B) returns C.
3 rule(dontPublishRaw).
4 has_target(dontPublishRaw, service15058189).
5 service(service15058189).
6 has_endpoint(service15058189,"http[s]?://.+").
7 receives_label(dontPublishRaw,raw).
8 has_decision(dontPublishRaw, dec).
9 has_effect(dec, drop).
10 has_obligation(dec,
11      log("Preventing data leak. ", message)).
```

## VII. PROTOTYPE EVALUATION

We implemented and evaluated a prototype of the LUCON policy framework to assess its application under real-world conditions.

### A. Implementation

As a platform for our implementation we chose the *Trusted IoT Connector* platform[4] – an open source platform based on the Karaf OSGi framework[5] that uses the Apache Camel message routing and mediation engine to forward messages between sensors and "applications" in form of Linux containers. While the Trusted Connector has been chosen for its security features, our implementation does not depend on it but would also be compatible with any other message router like Apache NiFi or Spring Integrations and other edge platforms like Eclipse Kura[6].

Apache Camel is a rule-based engine to route messages in the form of so-called *Exchange* objects according to Enterprise Integration Patterns (EIP). Due to its support for more than 240 protocol adapters, including HTTP, OPC-UA, MQTT, it is well-suited for IoT scenarios where data from different sources must be unified. We hook into the Camel engine by implementing an *interceptor* component that is called between each step in a message route and may

drop, forward or alter any Exchange object. The interceptor acts as the Policy Enforcement Point (PEP) and interacts with the other components of the LUCON framework which have been implemented as OSGi services. If the message is allowed to pass, the PEP simply puts it back into the processing engine. If the decision is to drop the message, the interceptor removes it from the message route and in case of an error, it exceptionally terminates the route, allowing a graceful exception handling. Any obligation that is possibly bound to the policy decision refers to OSGi services that the PEP will invoke. As OSGi services are dynamic and can spawn and terminate at any time, the set of supported obligations may vary at runtime and policy authors must consider that the execution of obligations may fail by stating an alternative effect in the `otherwise` element.

The Policy Decision Point (PDP) includes a tuProlog engine to load policies as Prolog theories and run queries against them. tuProlog is a Java-native lightweight Prolog implementation that has been chosen because of its small footprint of only 294 KB and especially because of its ability to map Prolog predicates to Java functions. An example of such mappings is shown in line 1 of Listing 2 where a Prolog predicate *regex* is defined by a call to the respective Java *regex* function to support querying for regular expressions, e.g. over service endpoint URLs.

In total, the size of the LUCON policy engine amounts to a 3.1 MB OSGi bundle that is loaded into the Karaf platform, automatically detects all Camel instances and hooks its interceptor into their message routes. The policy parser and code generator is not part of that engine in order to keep its footprint low. This means, policy authors will write policies in LUCON DSL in a separate IDE and load the compiled policies into the engine. The LUCON IDE has been implemented as an Eclipse "product" i.e. a standalone version of the Eclipse IDE that includes the code generator and various assistants for authoring the policy.

### B. Data Flow Awareness of Services

The most prevalent question for the integration of a policy framework is to which extend the existing IoT system must be aware of the framework and actively support it. LUCON requires only a single integration point, which is a hook into the message routing engine – realized as a Camel interceptor in our prototype. In addition LUCON does not require services that are able to handle message labels. We distinguish services by three classes of message handling capabilities: *agnostic*, *preserving*, and *active*.

*a) Agnostic Services:* Agnostic services are unaware of any data flow control mechanism. That is, when messages are sent into an agnostic service, all message labels will be lost and the data flow tracing will break. Most existing services will fall into this category.

Agnostic services are supported by LUCON's capability to state transformation functions as part of the policy. As long as transformation functions are specified for a service,

both runtime enforcement and static validation of routes will work as described above.

*b) Preserving:* Even if they are not aware of any data flow control mechanism, some services are able to preserve labels attached to messages. That is, when data is sent into the service and retrieved at a later time, previously attached labels will still be intact and data flow tracing is not interrupted. As example for such services are databases or file systems which persist message labels along with data records.

As long as preserving services do not perform any operation that would change labels, no transformation function needs to be stated in the policy. Data flow tracing will not break at runtime and labels will be transported across service calls. Also static route validation will work, as the service does not affect labels and thus remains irrelevant with respect to path explorations in the message route graph.

*c) Data Flow Aware Services:* Data flow aware services are able to actively modify message labels. While today, the vast majority of services is not data flow aware, an example of such services has been proposed in [25]. These services can modify message labels in a more complex way than could be expressed by transformation functions in the policy. The service in [25] for instance, modifies message labels according to an internal "taint logic" that cannot be written as a transformation function. As a consequence, static route validation with data flow aware services is only possible if the service's labeling semantics is available in the same logic representation as the policy.

In general, data flow awareness of services directly relates to the trust in that service. A non-aware service does not require a high level of trust, since it would not be able to alter labels in a malicious way. Data flow aware services, on the contrary, are able to modify labels and could interfere with data flow that way. Consequently, data-flow aware services require additional mechanisms for trust establishment, such as a remote attestation or certification.

### C. Performance Evaluation

The most critical metric for a policy engine is the time needed to evaluate a policy decision request. As each step in a message route requires a policy decision, the engine must not introduce unacceptable delays and must scale with an increasing number of services and rules. We evaluated how the runtime of the policy decision point for evaluating a decision request scales with the number of policies and services. While we consider a few dozens of rules to be a realistic size in most applications, we chose a test range of 1-5,000 rules and services. All rules were set up to match all services so that every decision request would require an evaluation of every single rule, which is the worst case. The tests were run against our prototype implementation which uses the Java-based tuProlog engine and does not include any runtime optimizations. As Figure 5 shows, the time for evaluating the decision requests scales linearly with
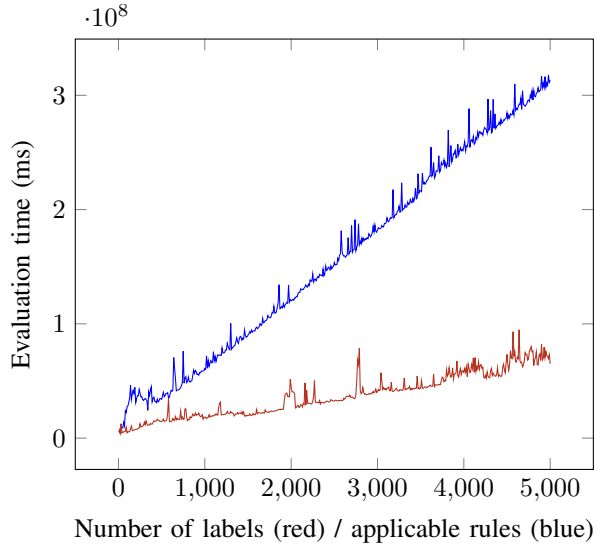


Figure 5. Policy decision time, scaling with rules (blue) and labels (red)

the number of rules within the analyzed range. For typical policy sizes of a few hundred rules, the evaluation takes approx. 12-15 ms. For a policy of 1,000 rules, it is still clearly below 50 ms and then increases linearly up to 150-200 ms for 5,000 rules. The red line in Figure 5 shows how runtime scales with an increasing number of message labels and a constant number of 50 rules. As can be seen, the decision time only depends on the number of rules, but does not increase with more labels.

The second metric of our performance evaluation is memory consumption. Here, we are especially interested if the framework is suited to run on typical IoT gateway devices or if the Prolog-based implementation it too memory-intensive for such applications. Figure 6 shows the memory consumption of the LUCON engine during a policy decision. Again, the blue line illustrates how memory consumption scales with an increasing number of rules and the red line indicates behavior with an increasing number of labels.

As expected, memory consumption scales linearly with the number of rules and constant with the number of labels, just as computation time does. The absolute numbers show that evaluating a policy of 50 rules requires less than 100 KB, while very large policies with thousands of rules may occupy several hundreds of megabyte of heap.

## VIII. CONCLUSIONS

In this paper we introduced LUCON, a policy framework for controlling data flows in distributed message-based systems. LUCON extends the concept of usage control by the notion of data flows. In contrast to traditional information flow control frameworks which enforce a single security model or information classification scheme, our approach labels messages and monitors their usage in a taint analysis style, addressing an attacker model in which information leaks via side channels such as observation
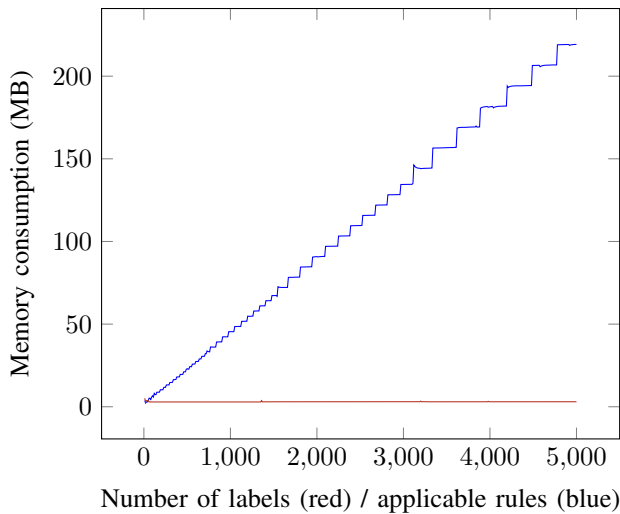
Figure 6. Memory for evaluating a policy decision. 1-5,00 rules/services/labels

of the control flow are negligible. An automated formal verification of LUCON policies against message routes informs users upfront about possible policy violations and thus supports policy authors in writing correct rules. Proofs created by the formal verification support system audits, as they assert that message routes will not violate security and privacy requirements.

Our prototype shows that the approach of compiling policies and message routes into the same logic representation is both suitable for runtime enforcement and static verification, without the need to convert back and forth between different representations and possible semantic gaps. A major question was if the performance of a Prolog-based evaluation engine can keep up with the demands of real-life systems with considerable high message throughput. Although performance impact of our prototype is notable, the measured delays in the range of 12-15 ms per policy decision are still in the range of typical network latency and suggest that with appropriate optimizations, the policy framework will easily be able to handle real world use cases.

### REFERENCES

[1] D. Basin, M. Harvan, F. Klaedtke, and E. Zălinescu. Monpoly: Monitoring usage-control policies. In *Proc. of the Second International Conference on Runtime Verification*, RV'11, pages 360–364, Berlin, Heidelberg, 2012. Springer-Verlag.

[2] D. Basin, F. Klaedtke, and S. Müller. Policy monitoring in first-order temporal logic. In *Computer Aided Verification*, volume 6174 of *Lecture Notes in Computer Science*, pages 1–18. Springer Berlin Heidelberg, 2010.

[3] D. E. Bell and L. J. LaPadula. Secure computer systems: Mathematical foundations. MITRE Corporation, 1973.

[4] K. J. Biba. Integrity considerations for secure computer systems. Technical report, MITRE Corp., 04 1977.

[5] G. Chinis, P. Pratikakis, S. Ioannidis, and E. Athanasopoulos. Practical information flow for legacy web applications. In *Proc. of the 8th Workshop on Implementation, Compilation, Optimization of Object-Oriented Languages, Programs and Systems*, pages 17–28. ACM, 2013.

[6] B. Davis and H. Chen. Dbtaint: cross-application information flow tracking via databases. *Proc. of WebApps*, 10, 2010.

[7] D. Denning. A lattice model of secure information flow. *Communications of the ACM*, 19(5):236–242, 1976.

[8] Y. Elrakaiby and J. Pang. Dynamic analysis of usage control policies. In *11th Int. Conf. on Security and Cryptography (SECRYPT)*, pages 88–100, Vienna, Austria, Nov. 2014.

[9] M. Harvan and A. Pretschner. State-based usage control enforcement with data flow tracking using system call interposition. In *Network and System Security, 2009. NSS '09. Third International Conference on*, pages 373–380, Oct 2009.

[10] M. Hilty, A. Pretschner, D. Basin, C. Schaefer, and T. Walter. A policy language for distributed usage control. In *ESORICS*, volume 4734, pages 531–546. Springer, 2007.

[11] G. Hohpe and B. Woolf. *Enterprise Integration Patterns: Designing, Building, and Deploying Messaging Solutions*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2003.

[12] B. Katt, X. Zhang, R. Breu, M. Hafner, and J.-P. Seifert. A general obligation model and continuity: enhanced policy enforcement engine for usage control. In *Proc. of the 13th ACM Symposium on Access Control Models and Technologies (SACMAT)*, pages 123–132. ACM, 2008.

[13] A. Lazouski, F. Martinelli, and P. Mori. Usage control in computer security: A survey. *Computer Science Review*, 4(2):81 – 99, 2010.

[14] A. C. Myers and B. Liskov. A decentralized model for information flow control. In *Proc. of the Sixteenth ACM Symposium on Operating Systems Principles*, SOSP '97, pages 129–142, New York, NY, USA, 1997. ACM.

[15] A. C. Myers, L. Zheng, S. Zdancewic, S. Chong, , and N. Nystrom. Jif: Java information flow. Software release, July 2001].

[16] J. Park and R. Sandhu. The $UCON_{ABC}$ usage control model. *ACM Trans. Inf. Syst. Secur.*, 7(1):128–174, Feb. 2004.

[17] T. Pasquier, J. Bacon, J. Singh, and D. Eyers. Data-centric access control for cloud computing. In *Proc. of the 21st ACM on Symposium on Access Control Models and Technologies*, SACMAT '16, pages 81–88, New York, NY, USA, 2016. ACM.

[18] T. F. J. Pasquier, J. Singh, D. M. Eyers, and J. Bacon. Camflow: Managed data-sharing for cloud services. *CoRR*, abs/1506.04391, 2015.

[19] T. F. J.-M. Pasquier, J. Bacon, and D. Eyers. FlowK: Information Flow Control for the Cloud. *6th Int. Conference on Cloud Computing Technology and Science (CloudCom)*, pages 1–8, 2014.

[20] A. Pretschner, M. Büchler, M. Harvan, C. Schaefer, and T. Walter. Usage control enforcement with data flow tracking for x11. In *Proc. of 5th Intl. Workshop on Security and Trust Management*, pages 124–137, 2009.

[21] A. Pretschner, M. Hilty, and D. Basin. Distributed usage control. *Communications of the ACM*, 49(9):39–44, 2006.

[22] A. Pretschner, J. Ruesch, C. Schaefer, and T. Walter. Formal analyses of usage control policies. In *Availability, Reliability and Security, 2009. ARES '09*, pages 98–105, March 2009.

[23] A. Sabelfeld and A. Russo. *From Dynamic to Static and Back: Riding the Roller Coaster of Information-Flow Control Research*, pages 352–365. Springer, Berlin, Heidelberg, 2010.

[24] R. Sandhu and J. Park. Usage control: A Vision for Next Generation Access Control. In *MMM-ACNS*, volume 2776, pages 17–31. Springer, 2003.

[25] J. Schütte and G. S. Brost. A data usage control system using dynamic taint tracking. In *Proc. of the Int. Conference on Advanced Information Network and Applications (AINA), year=2016, month = mar*.

[26] V. Simonet. The flow caml system. Software release, July 2003.

[27] X. Zhang, F. Parisi-Presicce, R. Sandhu, and J. Park. Formal model and policy specification of usage control. *ACM Transactions on Information and System Security (TISSEC)*, (4), Nov. 2005.

[28] X. Zhang, J.-P. Seifert, and R. Sandhu. Security enforcement model for distributed usage control. In *Sensor Networks, Ubiquitous and Trustworthy Computing (SUTC)*, pages 10–18, 2008.

Figure 7.   Operational semantics of dynamically taint-controlled message routes

$$\frac{\Delta, \cdot \vdash \texttt{get\_from(url)} \Downarrow \mu_m \qquad \tau'[m \leftarrow \mathcal{L}^+(url)] \qquad \iota' = \Sigma[pc+1]}{\langle \tau, \Sigma, \cdot, \Delta, \texttt{pc}, \iota \rangle, \texttt{from(url)} \rightarrow \langle \tau', \Sigma, \mu_m, \Delta, \texttt{pc+1}, \iota' \rangle} \text{ FROM}$$

$$\frac{\Delta, \mu_m \vdash \texttt{get\_from(url)} \Downarrow \mu_m \qquad \tau'[m \leftarrow \tau[m] \setminus \mathcal{L}^-(url) \cup \mathcal{L}^+(url)] \qquad \iota' = \Sigma[pc+1]}{\langle \tau, \Sigma, \mu_m, \Delta, \texttt{pc}, \iota \rangle, \texttt{to(url)} \rightarrow \langle \tau', \Sigma, \mu_m, \Delta, \texttt{pc+1}, \iota' \rangle} \text{ TO}$$

$$\frac{\Delta, \mu_m \vdash e \Downarrow 1 \qquad \Delta, \mu_m \vdash e_0 \Downarrow v_0 \qquad \iota' = \Sigma[v_0]}{\langle \tau, \Sigma, \mu_m, \Delta, \texttt{pc}, \iota \rangle, \textsf{when } e \textsf{ then goto } e_0 \textsf{ otherwise goto } e_1 \rightarrow \langle \tau, \Sigma, \mu_m, \Delta, v_0, \iota' \rangle} \text{ CHOICE (TRUE)}$$

$$\frac{\Delta, \mu_m \vdash e \Downarrow 0 \qquad \Delta, \mu_m \vdash e_1 \Downarrow v_1 \qquad \iota' = \Sigma[v_1]}{\langle \tau, \Sigma, \mu_m, \Delta, \texttt{pc}, \iota \rangle, \textsf{when } e \textsf{ then goto } e_0 \textsf{ otherwise goto } e_1 \rightarrow \langle \tau, \Sigma, \mu_m, \Delta, v_1, \iota' \rangle} \text{ CHOICE (FALSE)}$$

$$\frac{\Delta, \mu_m \vdash e \Downarrow m_0, ..., m_n \qquad \tau' = \tau \cup \bigcup_{0 \leq i \leq n} \tau'[i] = \tau[m] \qquad \iota' = \Sigma[pc+1]}{\langle \tau, \Sigma, \mu_m, \Delta, \texttt{pc}, \iota \rangle, \textsf{split } e \rightarrow \langle \tau', \Sigma, \mu_m, \Delta, pc+1, \iota' \rangle} \text{ SPLIT}$$

$$\frac{\Delta, \bigcup_i \mu_i, \vdash e \Downarrow m \qquad \tau'[m] = \bigcup_{0 \leq i \leq n} \tau[i] \qquad \iota' = \Sigma[pc+1]}{\langle \tau, \Sigma, \bigcup_i \mu_i, \Delta, \texttt{pc}, \iota \rangle, \textsf{aggregate } e \rightarrow \langle \tau', \Sigma, \mu_m, \Delta, pc+1, \iota' \rangle} \text{ AGGREGATE}$$

$$\frac{\Delta, \mu_m, \vdash e \Downarrow v \qquad \mu'_m[k \leftarrow v] \qquad \iota' = \Sigma[pc+1]}{\langle \tau, \Sigma, \bigcup_i \mu_i, \Delta, \texttt{pc}, \iota \rangle, \textsf{set-msg-prop}(k, e) \rightarrow \langle \tau, \Sigma, \mu'_m, \Delta, pc+1, \iota' \rangle} \text{ SET-MSG-PROP}$$

$$\frac{\Delta, \mu_m, \vdash e \Downarrow v \qquad \Delta'[k \leftarrow v] \qquad \iota' = \Sigma[pc+1]}{\langle \tau, \Sigma, \bigcup_i \mu_i, \Delta, \texttt{pc}, \iota \rangle, \textsf{set-env-prop}(k, e) \rightarrow \langle \tau, \Sigma, \mu_m, \Delta', pc+1, \iota' \rangle} \text{ SET-ENV-PROP}$$

$$\frac{\Delta, \mu_m \vdash \texttt{bean(b)} \Downarrow \mu_m \qquad \tau'[m \leftarrow \tau[m] \setminus \mathcal{L}^-(b) \cup \mathcal{L}^+(b)] \qquad \iota' = \Sigma[pc+1]}{\langle \tau, \Sigma, \mu_m, \Delta, \texttt{pc}, \iota \rangle, \texttt{bean(b)} \rightarrow \langle \tau', \Sigma, \mu_m, \Delta, \texttt{pc+1}, \iota' \rangle} \text{ BEAN}$$