



ProvMark: A Provenance Expressiveness Benchmarking System

Sheung Chi Chan
University of Edinburgh
United Kingdom
s1536869@inf.ed.ac.uk

James Cheney
University of Edinburgh
The Alan Turing Institute
United Kingdom
jcheney@inf.ed.ac.uk

Pramod Bhatotia
University of Edinburgh
The Alan Turing Institute
United Kingdom
pramod.bhatotia@ed.ac.uk

Thomas Pasquier
University of Bristol
United Kingdom
thomas.pasquier@bristol.ac.uk

Ashish Gehani
SRI International
United States
ashish.gehani@sri.com

Hassan Irshad
SRI International
United States
hassan.irshad@sri.com

Lucian Carata
University of Cambridge
United Kingdom
lc525@cam.ac.uk

Margo Seltzer
University of British Columbia
Canada
mseltzer@cs.ubc.ca

Abstract

System level provenance is of widespread interest for applications such as security enforcement and information protection. However, testing the correctness or completeness of provenance capture tools is challenging and currently done manually. In some cases there is not even a clear consensus about what behavior is correct. We present an automated tool, ProvMark, that uses an existing provenance system as a black box and reliably identifies the provenance graph structure recorded for a given activity, by a reduction to *subgraph isomorphism* problems handled by an external solver. ProvMark is a beginning step in the much needed area of testing and comparing the expressiveness of provenance systems. We demonstrate ProvMark's usefulness in comparing three capture systems with different architectures and distinct design philosophies.

CCS Concepts • **Applied computing** → **System forensics**; • **Software and its engineering** → **Operating systems**; *Software verification and validation*; • **Computing methodologies** → *Logic programming and answer set programming*.

Keywords provenance, benchmarking, subgraph isomorphism

ACM Reference Format:

Sheung Chi Chan, James Cheney, Pramod Bhatotia, Thomas Pasquier, Ashish Gehani, Hassan Irshad, Lucian Carata, and Margo Seltzer. 2019. ProvMark: A Provenance Expressiveness Benchmarking System. In *20th International Middleware Conference (Middleware '19)*, December 8–13, 2019, Davis, CA, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3361525.3361552>

1 Introduction

Data provenance is information about the origin, history, or derivation of some information [20]. It is commonly proposed as a basis

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 the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or to publish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Middleware '19, December 8–13, 2019, Davis, CA, USA

© 2019 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-7009-7/19/12...\$15.00

<https://doi.org/10.1145/3361525.3361552>

for reproducibility [28], dependability [3], and regulatory compliance [30], and it is increasingly being used for security, through forensic audit [33] or online dynamic detection of malicious behavior [16, 17]. To cater to the requirements of different use-cases, there are many system-level provenance capture systems in the literature, such as PASS [23], Hi-Fi [31], SPADE [12], OPUS [5], LPM [6], Inspector [32], and CamFlow [26] covering a variety of operating systems from Linux and BSD to Android and Windows. Each system assembles low level system events into a high-level *provenance graph* describing processes, system resources, and causal relationships among them.

Often, such systems are described as capturing a *complete* and *accurate* description of system activity. To date, the designers of each such system have decided how to interpret these goals independently, making different choices regarding what activity to record and how to represent it. Although some of these systems do use standards such as W3C PROV [7] that establish a common vocabulary for provenance-related data, such standards do *not* specify how to record operating system-level behaviour, or even when such records are considered “accurate” or “complete”. Indeed, as we shall see, in practice there is little consensus about how specific activities (e.g., renaming a file) should be represented in a provenance graph. Additionally, different systems also work at different system layers (e.g., kernel space vs. user space), so some information may be unavailable to a given system.

To illustrate the problem, consider Figure 1, which shows three provenance graph representations of the same rename system call, as recorded by three different systems. These graphs clearly illustrate nontrivial structural differences in how rename is represented as communications between processes (blue rectangles) and artifacts or resources (yellow ovals).

Security analysts, auditors, or regulators seeking to use these systems face many challenges. Some of these challenges stem from the inherent size and complexity of provenance graphs: it is not unusual for a single day's activity to generate graphs with millions of nodes and edges. However, if important information is missing, then these blind spots may render the records useless for their intended purpose. For example, if a provenance capture system does not record edges linking reads and writes to local sockets, then attackers can evade notice by using these communication

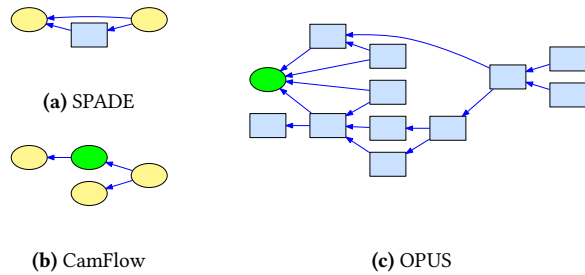


Figure 1. A rename system call, as recorded by three different provenance recorders. Images are clickable links to full-size versions with property labels.

channels. Such an omission could be due to a bug, but it could also result from misconfiguration, silent failure, or an inherent limitation of the recording system.

Provenance capture systems provide a broad-spectrum recording service, separate from the monitored applications, but (like any software system) they are not perfect, and can have their own bugs or idiosyncrasies. In order to rely on them for critical applications such as reproducible research, compliance monitoring, or intrusion detection, we need to be able to understand and validate their behavior. The strongest form of validation would consist of verifying that the provenance records produced by a system are accurate representations of the actual execution history of the system. However, while there is now some work on formalizing operating system kernels, such as seL4 [19] and HyperKernel [24], there are as yet no complete formal models of mainstream operating systems such as Linux. Developing such a model seems prerequisite to fully formalizing the accuracy, correctness, or completeness of provenance systems.

If there is no immediate prospect of being able to formally define and prove correctness for provenance recording systems, perhaps we can at least make it easier to compare and understand their behavior. We advocate a pragmatic approach as a first step toward the goal of validating and testing provenance systems, which (following Chan et al. [9]) we call *expressiveness benchmarking*. In contrast to performance benchmarking, which facilitates quantitative comparisons of run time or other aspects of a system’s behavior, expressiveness benchmarking facilitates qualitative comparisons of how different provenance systems record the same activities. The goal of such benchmarking is to answer questions such as:

- What does each node/edge in the graph tell us about actual events? (correctness, accuracy)
- What information is guaranteed to be captured and what are the blind spots? (completeness)

Previous work on comparing provenance representations has been based on manual inspection to compare the graphs [9, 21], but this manual approach is error prone and does not scale or allow automated testing. However, automating this task faces numerous challenges. It is difficult to configure some systems to record a single process’s activity reproducibly. Provenance records include volatile information such as timestamps and identifiers that vary across runs, even when the underlying process is deterministic. Furthermore, a process’s activity may include background process start-up activity that needs to be filtered out.

Expressiveness benchmarking is a *first (but important) step* towards understanding what it means for system provenance to be complete and/or correct. It offers a systematic means to compare different approaches. For example, how is a rename system call represented, and how do nodes and edges correspond to processes, files or file versions, or relationships among them? Are unsuccessful calls recorded? Are records of multiple similar calls aggregated together? By analyzing the same benchmark examples across different provenance capture systems, we can compare the information collected by different approaches to produce an understanding of the relative capabilities of the systems and their suitability for different tasks. Expressiveness benchmarking also enables *regression testing* since we can automatically compare benchmark graphs before and after a system change to detect differences.

We present *ProvMark*, an automated system for provenance expressiveness benchmarking. *ProvMark* executes a given target operation and records the resulting provenance captured by a given system. Different runs of the same system can then be compared to show whether each system records the target activity, and if so how. Our main contributions are as follows:

- We show how to generalize from multiple runs to obtain repeatable results, abstracting away volatile or transient data such as timestamps or identifiers.
- We solve the required graph matching problems using an external solver, based on Answer Set Programming (ASP), a variation of logic programming.
- We evaluate *ProvMark*’s performance and effectiveness in testing and comparing provenance systems, highlighting several bugs or idiosyncrasies found so far.

ProvMark has been developed in consultation with developers of three provenance recording systems, SPADE, OPUS and CamFlow, several of whom are coauthors of this paper. They have validated the results and in some cases helped adapt their systems to ease benchmarking. The *ProvMark* system along with supplementary results is publicly available at <http://provmark2018.github.io> [8].

2 Background and Related Work

Provenance [20] is a term originating in the art world. It refers to the chain of ownership of a work of art, including its original creator and each owner and location of the work from creation to the present. In the digital realm, provenance refers to information describing how an object came to be in its current form, typically including any other data from which it was derived and a description of the ways in which the input data was transformed to produce the output data. It is usually collected to allow post-hoc analysis to answer questions such as, “From where did this data come?”, “On what objects does it depend?”, or “What programs were used in its production?” Different applications of provenance will need different data, and systems are either tailored for specific applications or have a mechanism for selective capture, so that a user can obtain the right data for his/her intended use. Such flexibility makes it difficult to determine which system or configuration is most appropriate for a given task.

We apply *ProvMark* to three current provenance capture tools, namely SPADEv2 [12], OPUS [5] and CamFlow [26]. All three systems run under Linux. There are many other provenance capture systems. PASS is no longer maintained [23]; HiFi [31] was subsumed by LPM [6], which is similar to CamFlow, but less portable

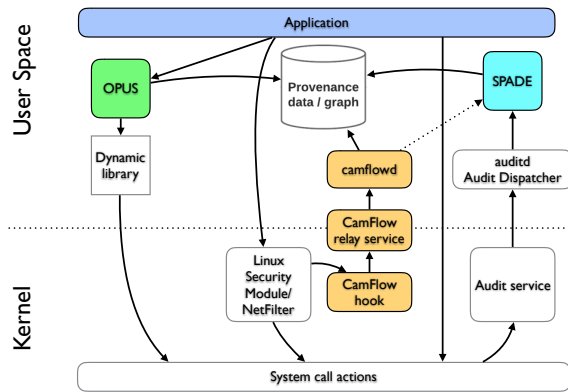


Figure 2. Architecture summary of the SPADE, OPUS, and CamFlow provenance recording systems

and not maintained. We focus on SPADE, OPUS and CamFlow as representative, currently available examples. Figure 2 shows how each system interacts with an application and the Linux kernel.

SPADE’s intended use is synthesizing provenance from machines in a distributed system, so it emphasizes relationships between processes and digital objects across distributed hosts. Our analysis uses SPADEv2 (tag *tc-e3*) with the Linux Audit Reporter [15], which constructs a provenance graph using information from the Linux audit system (including the Audit service, daemon, and dispatcher). SPADE runs primarily in user space and provides many alternative configurations, including filtering and transforming the data, for example, to enable versioning or finer-grained tracking of I/O or network events, or to make use of information in procsfs to obtain information about processes that were started before SPADE. We use a baseline configuration that collects data from the Audit Daemon (auditd), without versioning.

OPUS focuses on file system operations, attempting to abstract such operations and make the provenance collection process portable. It wraps standard C library calls with hooks that record provenance. The OPUS system is especially concerned with versioning support and proposes a Provenance Versioning Model, analogous to models previously introduced in the context of PASS [22] and later SPADE [13]. OPUS also runs in user space, but it relies on intercepting calls to a dynamically linked library (e.g., *libc*). Therefore, it is blind to activities that do not go through an intercepted dynamic library, but can observe the C library calls and userspace abstractions such as file descriptors.

CamFlow’s emphasis is sustainability through modularity, interfacing with the kernel via Linux Security Modules (LSM). The LSM hooks capture provenance, but then dispatch it to user space, via relays, for further processing. It strives for completeness and has its roots in Information Flow Control systems [29]. By default, CamFlow captures all system activity visible to LSM and relates different executions to form a single provenance graph; as we shall see, this leads to some complications for repeatable benchmarking. SPADE can also be configured to support similar behavior, while CamFlow can also be used (instead of Linux Audit) to report provenance to SPADE. Compared to SPADE and OPUS, which both run primarily in user space, CamFlow [26] monitors activity and generates the provenance graph from inside the kernel, via LSM and NetFilter hooks. This means the correctness of the provenance data

depends on the LSM operation. As the rules are set directly on the LSM hooks themselves, which are already intended to monitor all security-sensitive operations, CamFlow can monitor and/or record all sensitive operations. CamFlow allows users to set filtering rules when collecting provenance.

Prior work [9, 21] on comparing different provenance systems has followed a manual approach. For example, the Provenance Challenge [21] proposed a set of scientific computation scenarios and solicited submissions illustrating how different capture systems handle these scenarios, to facilitate (manual) comparison of systems. More recently, Chan et al. proposed a pragmatic, but also manual, approach [9] to benchmarking OS-level provenance tracking at the level of individual system calls. However, these manual approaches are error-prone and not scalable. It is also worth mentioning that Pasquier et al. [27] perform static analysis showing a conservative over-approximation of the CamFlow provenance recorded for each call. However, these results are not yet automatically compared with actual run-time behavior.

Provenance expressiveness benchmarking is also related to the emerging topic of *forensic-ready systems* [2, 25, 34]. Work in this area considers how to add logging to an existing system to detect known classes of behaviors, particularly for legal evidence or regulatory compliance, or how to advise developers on how to add appropriate logging to systems in development. The provenance collection systems above strive for completeness so that previously-unseen behaviors can be detected, and proactively record information so that user applications do not need to be modified.

3 System Design and Methodology

ProvMark is intended to automatically identify the (usually small) subgraph of a provenance graph that is recorded for a given target activity. Target activities could consist of individual system calls (or *syscalls*), sequences of syscalls, or more general (e.g., concurrent) processes. For the moment, we consider the simplest case of a single syscall, but the same techniques generalize to deterministic sequential target activities; handling concurrency and nondeterminism are beyond the scope of this paper.

We call the targeted system call the *target call* and the corresponding subgraph the *target graph*. Naively, one might proceed by writing a simple C program for each target system call that just performs that call and nothing else. However, starting and ending a process creates considerable “boilerplate” provenance, including calls such as *fork*, *execve*, and *exit*, as well as accesses to program files and libraries and, sometimes, memory mapping calls. Furthermore, some target calls require other *prerequisite* calls to be performed first. For example, analyzing a *read* or *close* system call requires first performing an *open*. Process startup and prerequisite calls are both examples of *background activity* that we would like to elide.

In each benchmark, we use an `#ifdef TARGET CPP` directive to identify the target behavior of interest. ProvMark generates two executables for each such benchmark: a *foreground program* that includes all code in the benchmark program, including the target and any context needed for it to execute, and a *background program* that contains the background activities. The two binaries are almost identical; the difference between the resulting graphs should precisely capture the target behavior.

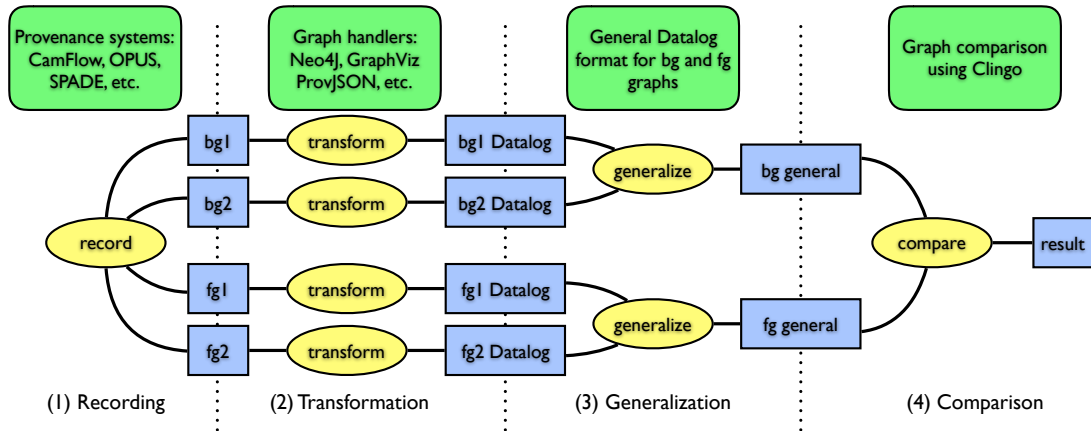


Figure 3. ProvMark system overview. The recording stage (1) uses one of the provenance recorders to compute background graphs bg_1, bg_2 and foreground graphs fg_1, fg_2 . (The same recorder is used for all four graphs.) The transformation stage (2) maps these graphs to a uniform Datalog format. The generalization stage (3) identifies the common structure of bg_1 and bg_2 , resulting in bg , and likewise fg_1 and fg_2 are generalized to fg . Finally, bg and fg are compared (4); structure corresponding to bg in fg is removed, yielding the benchmark result.

ProvMark includes a script for each syscall that generates and compiles the appropriate C executables and prepares a staging directory in which they will be executed with any needed setup, for example, first creating a file to run an `unlink` system call. We constructed these scripts manually since different system calls require different setup. The following code snippet illustrates the background program (including `open`) needed for the target `close` syscall (with `#ifdef` surrounding the target):

```
// close.c
#include <fcntl.h>
#include <unistd.h>
void main() {
    int id=open("test.txt", O_RDWR);
#ifdef TARGET
    close(id);
#endif
}
```

Figure 3 provides an overview of ProvMark, which is composed of four subsystems: (1) recording, (2) transformation, (3) generalization, and (4) comparison. Users can select which provenance capture system to use, which benchmark to run, and other configuration settings, such as the number of trials. Before presenting the details of the four subsystems, we outline some common use cases for ProvMark.

3.1 Use Cases

People are beginning to build critical distributed applications, particularly security applications, using system-level provenance. However, it is difficult for them to know how to interpret results or implement queries to detect activity on different systems. Both potential users and developers of all three considered provenance recording tools have participated in the design of ProvMark. To clarify when, and to whom, ProvMark is useful, in this section we outline several common use cases. In each case, ProvMark automates a central, labor-intensive step: namely, running tools to extract provenance, graphs and analyzing the graphs produced by a tool to identify a target activity. The first two use cases are real (albeit lightly fictionalized) situations where we have used ProvMark.

The other two are more speculative, but illustrate how ProvMark could be used for testing or exploration.

Tracking failed calls Alice, a security analyst, wants to know which provenance recorders track syscalls that fail due to access control violations, since these calls may be an indication of an attack or surveillance attempt. She can write small benchmark programs that capture various access control failure scenarios; most only take a few minutes to write, by modifying other, similar benchmarks for successful calls. Once this is done, ProvMark can run all of them to produce benchmark results. For example, Alice considers what happens if a non-privileged user unsuccessfully attempts to overwrite `/etc/passwd` by renaming another file.

By default SPADE installs Linux Audit rules that only report on successful system calls, so SPADE records no information in this case. OPUS monitors system calls via intercepting C library calls, so it knows whether a call is being attempted, and typically generates some graph structure even for unsuccessful calls. For example, the result of a failed `rename` call has the same structure as shown in Figure 1, but with a different return value property of `-1` instead of `0`. Finally, CamFlow can in principle monitor failed system calls, particularly involving permission checks, but does not do so in this case. Alice concludes that for her application, OPUS may provide the best solution, but resolves to discuss this issue with the SPADE and CamFlow developers as well.

Configuration validation Bob, a system administrator, wants to make sure that an installation of SPADE is configured correctly to match a security policy. SPADE (like the other systems) has several configuration parameters, to allow greater or lesser levels of detail, coalescing similar operations (such as repeated reads or writes), or enabling/disabling versioning. These also affect performance and system overhead, so Bob wants to ensure that enough information is recorded to enable successful audits, while minimizing system load.

Bob can use ProvMark to benchmark alternative configurations of SPADE. For example, SPADE provides a flag `simplify` that is

enabled by default. Disabling `simplify` causes `setresgid` and `setresuid` (among others) to be explicitly monitored, and Bob wants to ensure these calls are tracked. However, on doing so, Bob also uncovered a minor bug: when `simplify` is disabled, one of the properties of a background edge is initialized to a random value, which shows up in the benchmark as a disconnected subgraph. The bug was reported to the SPADE developers and quickly fixed.

SPADE also provides various *filters* and *transformers* which perform pre-processing or post-processing of the stored provenance respectively. Bob also experimented with one of SPADE’s filters, `IORuns`, which controls whether runs of similar read or write operations are coalesced into a single edge. Surprisingly, Bob found that in the benchmarked version of SPADE, enabling this filter had no effect. This turned out to be due to an inconsistency in the property names used by the filter vs. those generated by SPADE. This has also now been fixed.

Regression testing Charlie, a developer of provenance recording tool `XYZTrace`, wants to be able to document the level of completeness of `XYZTrace` to (potentially skeptical) users. `ProvMark` can be used for regression testing, by recording the graphs produced in a given benchmarking run, and comparing them with the results of future runs, using the same code for graph isomorphism testing `ProvMark` already uses during benchmarking. Charlie writes a simple script that stores the benchmark graphs (as `Datalog`) from previous runs, and whenever the system is changed, a new benchmarking run is performed and the results compared with the previous ones. When differences are detected, if the changes are expected then the new version of the graph replaces the old one; if the changes are unexpected, this is investigated as a potential bug.

Suspicious activity detection Dora, a security researcher, wants to identify patterns in provenance graphs that are indicative of a potential attack on the system. She compiles a set of scripts that carry out different kinds of attacks, and configures `CamFlow` on a virtual machine. She is particularly interested in detecting privilege escalation events where an attacker is able to gain access to new resources by subverting a privileged process. She instruments the scripts to treat the privilege escalation step as the “target activity”. Using `ProvMark`, Dora can obtain example provenance graphs that include the target activity or exclude it. If the graphs are not too large (e.g. hundreds rather than thousands of nodes), `ProvMark`’s default behavior will also be able to compute the differences between the graphs with and without the target activity.

3.2 Recording

The *recording* subsystem runs the provenance capture tools on test programs. This subsystem first prepares a staging directory that provides a consistent environment for test execution. The recording subsystem then starts the provenance capture tool with appropriate settings, captures the provenance generated by the tool, and stops the tool afterwards.

The recording subsystem is the only one to interact directly with the target provenance capture tools. For each tool, we implemented a script that configures the tool to capture the provenance of a specific process, rather than recording all contemporaneous system events. Recording multiple runs of the same process using `CamFlow` was challenging in earlier versions because `CamFlow` only serialized nodes and edges once, when first seen. The current version (0.4.5) provides a workaround to this problem that re-serializes the needed

Listing 1. Datalog Graph Format

```
Node n<gid>(<nodeID>, <label>)
Edge e<gid>(<edgeID>, <srcID>, <tgtID>, <label>)
Property p<gid>(<nodeID/edgeID>, <key>, <value>)
```

structures when they are referenced later. We also modified its configuration slightly to avoid tracking `ProvMark`’s own behavior. We otherwise use the default configuration for each system; we refer to these configurations as the *baseline* configurations.

The recording subsystem is used to record provenance graphs for the foreground and background variants of each benchmark. Provenance can include transient data that varies across runs, such as timestamps, so we typically record multiple runs of each program and filter out the transient information in a later *generalization* stage, described below.

Some of the provenance collecting tools were originally designed to record whole system provenance. The tools start recording when the machine is started and only stops when the machine is shut down. This behaviour ensures that provenance recording covers the entire operating system session. As we need to obtain repeatable results from several recording attempts for generalization and benchmark generation, we need to restart the recording section multiple times during the provenance collection. This may interfere with the results of some tools as they are not originally designed for frequent restarting of recording sessions. Thus the recording subsystem aims to manage the collecting tools by inserting and managing timeouts between the sessions. For example, we usually obtain stable results from SPADE by inserting timeouts to wait for successful completion of SPADE’s graph generation process; even so, we sometimes stop too early and obtain inconsistent results leading to mismatched graphs. Similarly, using `CamFlow`, we sometimes experience small variations in the size or structure of the results for reasons we have not been able to determine. In both cases we deal with this by running a larger number of trials and retaining the two smallest consistent results (as discussed later). For `OPUS`, any two runs are usually consistent, but starting `OPUS` and loading data into or out of `Neo4j` are time-consuming operations.

3.3 Transformation

Different systems output their provenance graphs in different formats. For example, SPADE supports `Graphviz DOT` format and `Neo4j` storage (among others), `OPUS` also supports `Neo4j` storage, and `CamFlow` supports `W3C PROV-JSON` [18] as well as a number of other storage or stream processing backends. `CamFlow` also can be used instead of `Linux Audit` as a reporter to SPADE, though we have not yet experimented with this configuration. To streamline the remaining stages, we translate these three formats to a common representation. Unlike the recording stage, this stage really is straightforward, but we will describe the common format in detail because it is important for understanding the remaining stages.

The common target format is a logical representation of property graphs, in which nodes and edges can have labels as well as associated properties (key-value dictionaries). Specifically, given a set Σ of node and edge labels, Γ of property keys, and D of data values, we consider property graphs $G = (V, E, src, tgt, lab, prop)$ where V is a set of vertex identifiers and E a set of edge identifiers; these are disjoint ($V \cap E = \emptyset$). Further, $src, tgt : E \rightarrow V$ maps each edge e to

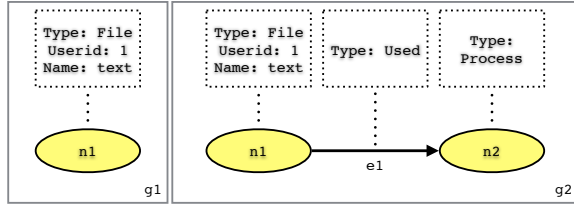


Figure 4. Sample Graphs

Listing 2. Datalog Format for Figure 4

```

ng1(n1, "File").
pg1(n1, "Userid", "1").
pg1(n1, "Name", "text").

ng2(n1, "File").
ng2(n2, "Process").
pg2(n1, "Userid", "1").
eg2(e1, n1, n2, "Used").
pg2(n1, "Name", "text").

```

its source and target nodes, respectively, $lab : V \cup E \rightarrow \Sigma$ maps each node or edge x to its label $lab(x) \in \Sigma$, and $prop : (V \cup E) \times \Gamma \rightarrow D$ is a partial function such that for a given node or edge x , then $prop(x, p)$ (if defined) is the value for property $p \in \Gamma$. In practice, Σ , Γ and D are each sets of strings.

For provenance graphs, following the W3C PROV vocabulary, the node labels are typically *entity*, *activity* and *agent*, edge labels are typically relations such as *wasGeneratedBy* and *used*, and properties are either PROV-specific property names or domain-specific ones, and their values are typically strings. However, our representation does not assume the labels and properties are known in advance; it works with those produced by the tested system.

We represent property graphs as sets of logical facts using a Prolog-like syntax called *Datalog* [1], which is often used to represent relational data in logic programming (as well as databases [1] and networking [14]).

The generic form of the Datalog graph format we use is shown as Listing 1. We assume a fixed string gid used to uniquely identify a given graph as part of its node, edge and label relation names. Each node $v \in V$ is represented as a fact $ngid(v, lab(v))$. Likewise, each edge $e = (v, w) \in E$ is represented as a fact $egid(e, src(e), tgt(e), lab(e))$. Finally, if a node or edge x has property p with value s , we represent this with the fact $pgid(x, p, s)$. Two sample graphs are shown in Figure 4 and their Datalog representations in Listing 2.

All the remaining stages, including the final result, work on the Datalog graph representation, so these stages are independent of the provenance capture tool and its output format. The Datalog representation can easily be visualized.

3.4 Graph Generalization

The third subsystem performs *graph generalization*. Recall that the recording stage produces several graphs for a given test program. We wish to identify a single, representative and general graph for each test program. To formalize these notions, we adapt the notion of *graph isomorphism* to property graphs: G_1 is isomorphic to G_2 if there is an invertible function $h : V_1 \cup E_1 \rightarrow V_2 \cup E_2$ such that

1. $h(src_1(e)) = src_2(h(e))$ and $h(tgt_1(e)) = tgt_2(h(e))$, i.e. h preserves edge relationships

2. $lab_1(x) = lab_2(h(x))$, i.e. h preserves node and edge labels
3. $prop_1(x, k, v) = prop_2(h(x), k, v)$, i.e. h preserves properties.

Moreover, we say that G_1 and G_2 are *similar* if only the first two conditions hold, that is, if G_1 and G_2 have the same shape but possibly different properties.

We assume that over sufficiently many recording trials, there will be at least two *representative* ones, which are similar to each other. Identifying two representative runs is complicated by the fact that there might be multiple pairs of similar graphs. We discuss a strategy for identifying an appropriate pair below. To obtain two representative graphs, we first consider all of the trial runs and partition them into similarity classes. We first discard all graphs that are only similar to themselves, and consider these to be failed runs. Among the remaining similarity classes, we choose a pair of graphs whose size is smallest. Picking the two largest graphs also seems to work; the choice seems arbitrary. However, picking the largest background graph and the smallest foreground graph leads to failure if the extra background structure is not found in the foreground, while making the opposite choice leads to extra structure being found in the difference.

Given two similar graphs, the generalization stage identifies the property values that are consistent across the two graphs and removes the transient ones. The generalization stage searches for a matching of nodes and edges of the two graphs that minimizes the number of different properties. We assume that the differences are transient data and discard them. We would like to minimize the number of differences, that is, match as many properties as possible.

Similarity and generalization are instances of the *graph isomorphism* problem, whose status (in P or NP-complete) is unknown [4]. We solve these problems using an Answer Set Programming (ASP) [11] specification that defines the desired matching between the two graphs. ASP is a decidable form of logic programming that combines a high-level logical specification language with efficient search and optimization techniques analogous to SAT solving or integer linear programming. An answer set program specifies a set of possible *models*, which correspond to solutions to a problem. The specification is a collection of logical formulas that define when a model is a possible solution. In our setting, the models consist of matching relations between two graphs, satisfying the requirements of graph similarity. ASP is well-suited to this problem because the graphs and specification can be represented concisely, and (as we shall see) the ASP solver can efficiently find solutions.

The problem specification is a logic program defining when a binary relation forms an isomorphism between the graphs. The code in Listing 3 defines graph isomorphism in ASP. ASP specifications consist of rules that constrain the possible solution set. The first four lines specify that h can relate any node in G_1 to any node in G_2 and vice versa, and similarly for edges. The next two lines constrain h to be a 1-1 function. (Rules of the form $:- A, B, C.$ can be read as "It cannot be the case that A, B, and C hold".) The remaining three pairs of lines specify that h must preserve node and edge labels and the source and targets of edges.

This specification can be solved using *clingo*, an efficient ASP solver [11]. As ASP is a kind of logic programming, it helps ProverMark to determine an optimal matching solution among two graphs by searching for a model that satisfies the specification. We use the resulting matching to determine which properties are common

Listing 3. Graph similarity

```

{h(X,Y) : n2(Y,_)} = 1 :- n1(X,_).
{h(X,Y) : n1(X,_)} = 1 :- n2(Y,_).
{h(X,Y) : e2(Y,_,_,_)} = 1 :- e1(X,_,_,_).
{h(X,Y) : e1(X,_,_,_)} = 1 :- e2(Y,_,_,_).

:- X <> Y, h(X,Z), h(Y,Z).
:- X <> Y, h(Z,Y), h(Z,X).

:- n1(X,L), h(X,Y), not n2(Y,L).
:- n2(Y,L), h(X,Y), not n1(X,L).

:- e1(E1,_,_,L), h(E1,E2), not e2(E2,_,_,L).
:- e2(E2,_,_,L), h(E1,E2), not e1(E1,_,_,L).

:- e1(E1,X,_,_), h(E1,E2), e2(E2,Y,_,_), not h(X,Y).
:- e1(E1,_,X,_), h(E1,E2), e2(E2,_,Y,_), not h(X,Y).

```

across both graphs and discard the others. We perform generalization independently on the foreground and background graphs. The outputs of this stage are the two generalized graphs representing the invariant activity of the foreground and background programs respectively.

3.5 Graph Comparison

The fourth and last subsystem is *graph comparison*. Its purpose is to match the background graph to a subgraph of the foreground graph; the unmatched part of the foreground graph corresponds to the target activity.

Because provenance recording is generally monotonic (append-only), we expect that the generalized provenance graph for the background program is a subgraph of the generalized provenance graph for the foreground program, so there should be a one-to-one matching from the nodes and edges in the background graph to the foreground graph. This graph matching problem is a variant of the subgraph isomorphism problem, a classical NP-complete problem [10]. We again solve these problems using ASP, taking advantage of the fact that ASP solvers can search for optimal solutions according to some cost model.

Given two graphs, G_1 and G_2 , the solver finds a matching from the nodes and edges of G_1 to those of G_2 that identifies a subgraph of G_2 isomorphic to G_1 and minimizes the number of mismatched properties. Listing 4 shows the approximate subgraph optimization problem specification. It is related to graph similarity, but only requires nodes/edges in G_1 to be matched to one in G_2 and not the reverse. Additionally, the last four lines define the cost of a matching as the number of properties present in G_1 with no matching property in G_2 , and this cost is to be minimized.

We again use the clingo solver to find an optimal matching. Once we have found such a matching, the graph comparison subsystem subtracts the matched nodes and edges in the background activity from the foreground graph. This is essentially a set difference operation between the nodes and edges of the graphs, but we also retain any nodes that are sources or targets of edges in the difference; these are visualized as green or gray *dummy nodes*.

Listing 4. Approximate subgraph isomorphism

```

{h(X,Y) : n2(Y,_)} = 1 :- n1(X,_).
{h(X,Y) : e2(Y,_,_,_)} = 1 :- e1(X,_,_,_).

:- X <> Y, h(X,Z), h(Y,Z).
:- X <> Y, h(Z,Y), h(Z,X).

:- n1(X,L), h(X,Y), not n2(Y,L).
:- e1(E1,_,_,L), h(E1,E2), not e2(E2,_,_,L).

:- e1(E1,X,_,_), h(E1,E2), e2(E2,Y,_,_), not h(X,Y).
:- e1(E1,_,X,_), h(E1,E2), e2(E2,_,Y,_), not h(X,Y).

cost(X,K,0) :- p1(X,K,V), h(X,Y), p2(Y,K,V).
cost(X,K,1) :- p1(X,K,V), h(X,Y), p2(Y,K,W), V <> W.
cost(X,K,1) :- p1(X,K,V), h(X,Y), not p2(Y,K,_) .

#minimize { PC,X,K : cost(X,K,PC) }.

```

1	Files	close, creat, dup[2,3], [sym]link[at], mknod[at], open[at], [p]read, rename[at], [f]truncate, unlink[at], [p]write
2	Processes	clone, execve, exit, [v]fork, kill
3	Permissions	[f]chmod[at], [f]chown[at], set[re[s]]gid, set[re[s]]uid
4	Pipes	pipe[2], tee

Table 1. Benchmarked syscalls

4 Demonstration

We present a demonstration of using ProvMark to investigate and compare the behavior of the different capture systems. These results pertain to SPADEv2 (tag *tc-e3*) and OPUS version 0.1.0.26 running on Ubuntu 16.04 and CamFlow version 0.4.5 running on Fedora 27.

Unix-like operating systems support over three hundred system calls. We focus on a subset of 22 common system call families shown in Table 1, including the most common file and process management calls. We prepared benchmarking programs for these syscalls, along with tests for each one to ensure that the target behavior was performed successfully. Note that the same system call may display different behavior using different parameters and system states; we focus on common-case behavior here, but handling other scenarios such as failure cases is straightforward as outlined in Section 3.1.

We cannot show all of the benchmark graphs, nor can we show the graphs at a readable size; the full results are available for inspection online [8]. Table 2 summarizes the results. In this table, “ok” means the graph is correct (according to the system’s developers), and “empty” means the foreground and background graphs were similar and so the target activity was undetected. The table includes notes indicating the reason for emptiness. Also, Table 3 shows selected graphs (again, images are links to scalable online images with full details) for the benchmark result. Although the property keys and values are omitted, we can compare the graph structure and see that the three tools generate very different structures for the same system calls. This demonstrates one of the important motivations for expressiveness benchmarking. In these graph results, the yellow ovals represent artifacts which includes files, pipes or

other resources. The blue rectangles represent processes involved in the named system calls. The green or gray ovals represent dummy nodes which stand for pre-existing parts of the graph. We retain these dummy nodes to make the result a complete graph.

In this section we present and discuss the benchmarking outputs for several representative syscalls, highlighting minor bugs found and cases where the behavior of individual tools or the variation among different tools was surprising.

4.1 File system

File system calls such as `creat`, `open`, `close` are generally covered well by all of the tools, but they each take somewhat different approaches to representing even simple behavior. For example, for the `open` call, SPADE adds a single node and edge for the new file, while CamFlow creates a node for the file object, a node for its path, and several edges linking them to each other and the opening process, and OPUS creates four new nodes including two nodes corresponding to the file. On the other hand, reads and writes appear similar in both SPADE and CamFlow, while by default, OPUS does not record file reads or writes. For `close`, CamFlow records the underlying kernel data structures eventually being freed, but ProvMark does not reliably observe this.

The `dup[2,3]` syscall duplicates a file descriptor so that the process can use two different file descriptors to access the same open file. SPADE and CamFlow do not record this call directly, but the changes to the process's file descriptor state can affect future file system calls that are monitored. OPUS does record `dup[2,3]`. The two added nodes are not directly connected to each other, but are connected to the same process node in the underlying graph. One component records the system call itself and the other one records the creation of a new resource (the duplicated file descriptor).

The `link[at]` call is recorded by all three systems but `symlink[at]` is not recorded by CamFlow 0.4.5. Likewise, `mknod` is only recorded by OPUS, and `mknodat` is not recorded by any of the systems.

The `rename[at]` call, as shown in the introduction, illustrates how the three systems record different graph structure for this operation. SPADE represents a rename using two nodes for the new and old filenames, with edges linking them to each other and to the process that performed the rename. OPUS creates around a dozen nodes, including nodes corresponding to the rename call itself and the new and old filename. CamFlow represents a rename as adding a new path associated with the file object; the old path does not appear in the benchmark result.

4.2 Process management

The process management operations start processes (`[v]fork`, `clone`), replace a process's memory space with a new executable and execute it (`execve`), or terminate a process (`kill`, `exit`).

Process creation and initialization calls are generally monitored by all three tools, except that OPUS does not appear to monitor `clone`. Inspecting the results manually shows significant differences, however. For example, the SPADE benchmark graph for `execve` is large, but has just a few nodes for OPUS and CamFlow. On the other hand, the `fork` and `vfork` graphs are small for SPADE and CamFlow and large for OPUS. This may indicate the different perspectives of these tools: SPADE relies on the activity reported by Linux Audit, while OPUS relies on intercepting C library calls outside the kernel.

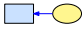
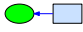


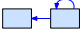
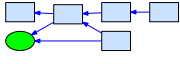
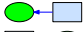
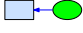


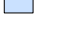


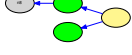





Group	syscall	SPADE	OPUS	CamFlow
1	close	ok	ok	empty (LP)
1	creat	ok	ok	ok
1	dup	empty (SC)	ok	empty (NR)
1	dup2	empty (SC)	ok	empty (NR)
1	dup3	empty (SC)	ok	empty (NR)
1	link	ok	ok	ok
1	linkat	ok	ok	ok
1	symlink	ok	ok	empty (NR)
1	symlinkat	ok	ok	empty (NR)
1	mknod	empty (NR)	ok	empty (NR)
1	mknodat	empty (NR)	empty (NR)	empty (NR)
1	open	ok	ok	ok
1	openat	ok	ok	ok
1	read	ok	empty (NR)	ok
1	pread	ok	empty (NR)	ok
1	rename	ok	ok	ok
1	renameat	ok	ok	ok
1	truncate	ok	ok	ok
1	ftruncate	ok	ok	ok
1	unlink	ok	ok	ok
1	unlinkat	ok	ok	ok
1	write	ok	empty (NR)	ok
1	pwrite	ok	empty (NR)	ok
2	clone	ok	empty (NR)	ok
2	execve	ok	ok	ok
2	exit	empty (LP)	empty (LP)	empty (LP)
2	fork	ok	ok	ok
2	kill	empty (LP)	empty (LP)	empty (LP)
2	vfork	ok (DV)	ok	ok
3	chmod	ok	ok	ok
3	fchmod	ok	empty (NR)	ok
3	fchmodat	ok	ok	ok
3	chown	empty (NR)	ok	ok
3	fchown	empty (NR)	empty (NR)	ok
3	fchownat	empty (NR)	ok	ok
3	setgid	ok	ok	ok
3	setregid	ok	ok	ok
3	setresgid	empty (SC)	empty (NR)	ok
3	setuid	ok	ok	ok
3	setreuid	ok	ok	ok
3	setresuid	ok (SC)	empty (NR)	ok
4	pipe	empty (NR)	ok	empty (NR)
4	pipe2	empty (NR)	ok	empty (NR)
4	tee	empty (NR)	empty (NR)	ok

Note	Meaning
NR	Behavior not recorded (by default configuration)
SC	Only state changes monitored
LP	Limitation in ProvMark
DV	Disconnected vforked process

Table 2. Summary of validation results

Another interesting observation is that the SPADE benchmark results for `fork` and `vfork` differ. Specifically, for `vfork`, SPADE represents the forked process as a disconnected activity node, i.e. there is no graph structure connecting the parent and child process. This is because Linux Audit reports system calls on `exit`, while the parent process that called `vfork` is suspended until the child

Table 3. Example benchmark results for SPADE, OPUS and CamFlow. Images are links to scalable online images with property labels.

	open	read	write	dup	setuid	setresuid
SPADE				Empty		
OPUS		Empty	Empty	 	  	Empty
CamFlow			 	Empty		  

process exits. SPADE sees the vforked child process in the log executing its own syscalls before it actually sees the vfork that created the child process.

The exit and kill calls are not detected because they deviate from the assumptions our approach is based on. A process always has an implicit exit at the end, while killing a process means that it does not terminate normally. Thus, the exit and kill benchmarks are all empty. We plan to consider a more sophisticated approach that can benchmark these calls in future work.

4.3 Permissions

We included syscalls that change file permissions (such as `[f]chown[at]`, `[f]chmod[at]`) or process permissions in this group. According to its documentation, SPADE currently records `[f]chmod[at]` but not `[f]chown[at]`. OPUS does not monitor `fchmod` or `fchown` because from its perspective these calls only perform read/write activity and do not affect the process’s file descriptor state, so as for other read/write activity OPUS does not record anything in its default configuration. CamFlow records all of these operations.

In its default configuration, SPADE does not explicitly record `setresuid` or `setresgid`, but it does monitor changes to these process attributes and records observed changes. Our benchmark result for `setresuid` is nonempty, reflecting an actual change of user id, while our benchmark for `setresgid` just sets the group id attribute to its current value, and so this activity is not noticed by SPADE. OPUS also does not track these two calls, while CamFlow again tracks all of them.

4.4 Pipes

Pipes provide efficient interprocess communication. The `pipe[2, 3]` call creates a pipe, which can be read or written using standard file system calls, and `tee` duplicates information from one pipe into another without consuming it. Of the three systems, only OPUS records `pipe[2, 3]` calls, while only CamFlow currently records `tee`.

5 System Evaluation

In this section, we will evaluate ProvMark. We claim that the ProvMark system makes it easy to test and compare the behavior of provenance tracking systems and to assess their completeness and correctness. The previous section gave some evidence that the results are useful for understanding the different systems and what information they do or do not collect. In this section, we will evaluate the performance and extensibility of ProvMark itself. For a fair comparison, all the experiments are done in a virtual machine

with 1 virtual CPU and 4GB of virtual memory. Dependencies for all three provenance collecting tools, including Neo4J and some needed Python libraries, are also installed. All virtual machines were hosted on a Intel i5-4570 3.2GHz with 8GB RAM running Scientific Linux 7.5.

5.1 System Performance

We first discuss the performance of the ProvMark system. Some performance overhead is acceptable for a benchmarking or testing tool; however, we need to establish that ProvMark’s strategy for solving NP-complete subgraph isomorphism problems is effective in practice (i.e. takes minutes rather than days). In this section, we will show that the extra work done by ProvMark to generate benchmarks from the original provenance result is acceptable.

We first report measurements of the recording time. Since the number of trials varies, we report the average time needed per trial. For SPADE, recording took approximately 20s per trial. For OPUS, the recording time was around 28s per trial, and for CamFlow each trial took around 10s. We wish to emphasize that the recording time results are *not* representative of the recording times of these systems in normal operation — they include repeatedly starting, stopping, and waiting for output to be flushed, and the waiting times are intentionally conservative to avoid garbled results. No conclusions about the relative performance of the recording tools when used as intended should be drawn.

In Figures 5–7, we summarize the time needed for ProvMark to run five representative syscalls using SPADE, OPUS, and CamFlow respectively. The bars are divided into three portions, representing the time needed for the transformation, generalization and comparison subsystems. Note that the x-axes are not all to the same scale: in particular the transformation, generalization and comparison times for OPUS are much higher than for the other tools because of database startup and access time, and because the graphs extracted from the database are larger. Modifying OPUS to circumvent the database and serialize provenance directly would avoid this cost, but we have chosen to use the tools as they are to the greatest extent possible.

From the data in Figures 5–7 we can see that the transformation stage is typically the most time-consuming part. The transformation stage maps the different provenance output formats to Datalog format. The transformation time for OPUS is much longer than for the other two systems. This appears to be because OPUS produces larger graphs (including recording environment variables), and extracting them involves running Neo4j queries, which also has a one-time JVM warmup and database initialization cost.

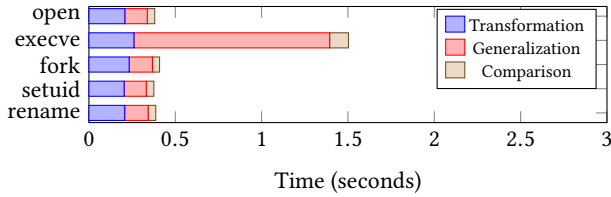


Figure 5. Timing results: SPADE+Graphviz

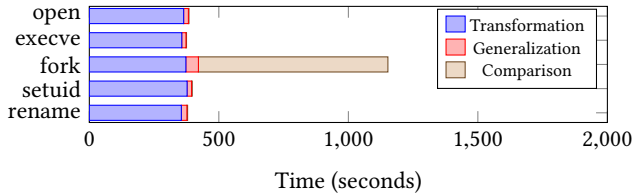


Figure 6. Timing results: OPUS+Neo4J

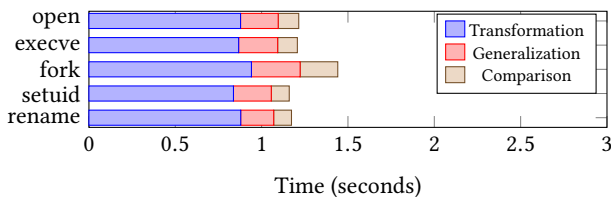


Figure 7. Timing results: CamFlow+ProvJson

The time required for the generalization stage depends mainly on the number of elements (nodes and edges) in the graph, since this stage matches elements in pairs of graphs to find an isomorphic graph matching. As we can see from the results, the generalization of OPUS graphs again takes significantly longer than for SPADE and CamFlow, probably because the graphs are larger. For SPADE, generalization of the `execve` benchmark takes much longer than for other calls (though still only a few seconds).

The time required for the comparison stage is usually less than for generalization (recall that in generalization we process two pairs of graphs whereas in comparison we just compare the background and foreground graphs). Comparison of OPUS graphs again takes longest, while comparison of CamFlow graphs again takes slightly longer than for SPADE, perhaps because of the larger number of properties.

In general, we can conclude that the time needed for ProvMark's transformation, generalization, and comparison stages is acceptable compared with the time needed for recording. Most benchmarks complete within a few minutes at most, though some that produce larger target graphs take considerably longer; thus, at this point running all of the benchmarks takes a few hours. This seems like an acceptable price to pay for increased confidence in these systems, and while there is clearly room for improvement, this compares favorably with manual analysis, which is tedious and would take a skilled user several hours. In addition, we have monitored the memory usage and found that ProvMark never used more than 75% of memory on a 4GB virtual machine, indicating memory was not a bottleneck.

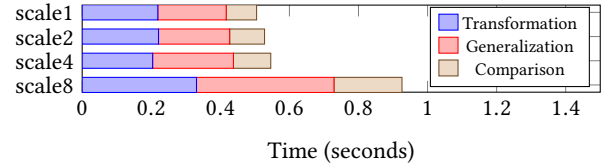


Figure 8. Scalability results: SPADE+Graphviz

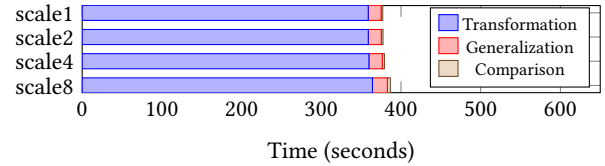


Figure 9. Scalability results: OPUS+Neo4J

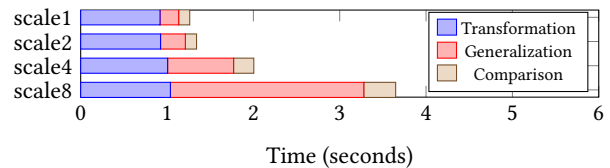


Figure 10. Scalability results: CamFlow+ProvJson

5.2 Scalability

Our experiments so far concentrated on benchmarking one syscall at a time. The design of ProvMark allows arbitrary activity as the target action, including sequences of syscalls. As mentioned above, we can simply identify a target action sequence using `#ifdef TARGET` in order to let ProvMark generate background and foreground programs respectively.

Of course, if the target activity consists of multiple syscalls, the time needed for ProvMark to process results and solve the resulting subgraph isomorphism problems will surely increase. The generated provenance graph will also increase in size and number of elements. The NP-complete complexity of sub-graph isomorphism means that in the worst case, the time needed to solve larger problem instances could increase exponentially. This section investigates the scalability of ProvMark when the size of the target action increases.

In figure 8-10, the charts show the time needed for the three processing subsystems of ProvMark in handling some scalability test cases on SPADE, OPUS and CamFlow respectively. The scalability test cases range from `scale1` to `scale8`. In test case `scale1`, the target action sequence is simply a creation of a file and another deletion of the newly created file. In test case `scale2`, `scale4` and `scale8`, the same target action is repeated twice, four times, and eight times respectively. The results show that the time needed initially grows slowly for SPADE, but by `scale8` the time almost doubles compared to `scale1`. For OPUS, the time increases are dwarfed by the high overhead of transformation, which includes the one-time Neo4j startup and access costs as discussed above. For CamFlow, the time needed approximately doubles at each scale factor. Although further scalability experiments are needed to consider much larger graph or benchmark sizes, these experiments do demonstrate that ProvMark can currently handle short sequences of 10-20 syscalls without problems.

Module (Format)	SPADE (DOT)	OPUS (Neo4j)	CamFlow (PROV-JSON)
Recording	171	118	192
Transformation	74	122	128

Table 4. Module sizes (Python lines of code)

5.3 Modularity and Extensibility

ProvMark was designed with extensibility in mind. Only the first two stages (recording and generalization) depend on the details of the recording system being benchmarked, so to support a new system, it suffices to implement an additional recording module that uses the new system to record provenance for a test executable, and (if needed) an additional translation module that maps its results to Datalog.

As discussed earlier, it was non-trivial to develop recording modules for the three systems that produce reliable results. In particular, supporting CamFlow required several iterations and discussion with its developers, which have resulted in changes to CamFlow to accommodate the needs of benchmarking. This architecture has been refined through experience with multiple versions of the SPADE and CamFlow systems, and we have generally found the changes needed to ProvMark to maintain compatibility with new versions to be minor. Developing new recording modules ought to be straightforward for systems that work in a similar way to one of the three currently supported.

If a provenance tool generates data in a format not currently supported by ProvMark, an additional module for handling this type of provenance data is needed. The need for new transformations should decrease over time as more tools are introduced to ProvMark since there are a limited number of common provenance formats.

As shown in Table 4, none of the three recording or transformation modules required more than 200 lines of code (Python 3 using only standard library imports). We initially developed ProvMark with support for SPADE/DOT and OPUS/Neo4j combinations; thus, adding support for CamFlow and PROV-JSON only required around 330 lines of code.

Finally, throughout the development period of ProvMark, the candidate tools have all been updated and include new features and behaviour. Over time, we have updated ProvMark to react to changes in both SPADE and CamFlow, with few problems.

5.4 Discussion and Limitations

While we have argued that ProvMark is effective, fast enough to produce timely results, and easy to extend, it also has some limitations, which we now discuss.

At this stage, ProvMark has only been applied to provenance tracking at the operating-system level; in particular, for SPADE, OPUS and CamFlow running under Linux. We believe the same approach can be adapted to other operating systems or other layers of distributed systems, assuming an extension to deal with nondeterministic or concurrent activity as described below.

We have evaluated ProvMark using a small subset of individual syscalls. Creating benchmarks currently takes some manual effort and it would be helpful to automate this process to ease the task of checking all of the possible control flow paths each syscall can take, or even generate multi-step sequences. In addition, the analysis and comparison of the resulting benchmark graphs requires some manual effort and understanding.

We have shown ProvMark is capable of handling benchmark programs with a small number of target system calls. Scaling up to larger amounts of target activity or background activity appears challenging, due to the NP-completeness and worst-case exponential behavior of subgraph isomorphism testing. However, for deterministic activity it may be possible to take advantage of other structural characteristics of the provenance graphs to speed up matching: for example, if matched nodes are usually produced in the same order (according to timestamps), then it may be possible to incrementally match the foreground and background graphs.

Lastly, ProvMark currently handles deterministic activity only. Nondeterminism (for example through concurrency) introduces additional challenges: both the foreground and background programs might have several graph structures corresponding to different schedules, and there may be a large number of different possible schedules. It also appears necessary to perform some kind of fingerprinting or graph structure summarization to group the different possible graphs according to schedule. We may also need to run larger numbers of trials and develop new ways to align the different structures so as to obtain reliable results. It appears challenging to ensure completeness, that is, that all of the possible behaviors are observed, especially since the number of possible schedules can grow exponentially in the size of the program.

6 Conclusions and Future Work

This paper presents ProvMark, an automated approach to benchmarking and testing the behavior of provenance capture systems. To the best of our knowledge it is the first work to address the unique challenges of testing and comparing provenance systems. We have outlined the design of ProvMark, and showed how it helps identify gaps in coverage, idiosyncrasies and even a few bugs in three different provenance capture systems. We also showed that it has acceptable performance despite the need to solve multiple NP-complete graph matching subproblems. ProvMark is a significant step towards validation of such systems and should be useful for developing correctness or completeness criteria for them.

There are several directions for future work, to address the limitations discussed in Section 5.4. First, additional support for automating the process of creating new benchmarks, or understanding and interpreting the results, would increase the usefulness of the system, as would extending it to provenance tracking at other levels, such as distributed system coordination layers. Second, although we have evaluated scalability up to 10–20 target system calls, realistic applications such as suspicious behavior analysis would require dealing with much larger target activities and graphs. Finally, ProvMark cannot deal with nondeterminism, including concurrent I/O and network activity, and extending its approach to handle nondeterministic target activity is a focus of current work.

Acknowledgments

Effort sponsored by the Air Force Office of Scientific Research, Air Force Material Command, USAF, under grant number FA8655-13-1-3006. The U.S. Government and University of Edinburgh are authorised to reproduce and distribute reprints for their purposes notwithstanding any copyright notation thereon. Cheney was also supported by ERC Consolidator Grant Skye (grant number 682315).

This material is based upon work supported by the Defense Advanced Research Projects Agency (DARPA) under contract FA8650-15-C-7557 and the National Science Foundation under Grant ACI-1547467 and SSI-1450277 (End-to-End Provenance). Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

We are grateful to Ajitha Rajan for comments on a draft of this paper and to our shepherd Boris Koldehofe and anonymous reviewers for insightful comments.

References

- [1] Serge Abiteboul, Richard Hull, and Victor Vianu. 1995. *Foundations of databases*. Addison-Wesley.
- [2] Dalal Alrajeh, Liliana Pasquale, and Bashar Nuseibeh. 2017. On evidence preservation requirements for forensic-ready systems. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering (ESEC/FSE 2017)*. 559–569.
- [3] Peter Alvaro and Severine Tymon. 2017. Abstracting the geniuses away from failure testing. *Queue* 15, 5 (2017), 10.
- [4] Vikraman Arvind and Jacobo Torán. 2005. Isomorphism testing: Perspectives and open problems. *Bulletin of the European Association for Theoretical Computer Science* 86 (2005), 66–84.
- [5] Nikilesh Balakrishnan, Thomas Bytheway, Ripduman Sohan, and Andy Hopper. 2013. OPUS: A Lightweight System for Observational Provenance in User Space. In *Proceedings of the 5th USENIX Workshop on Theory and Practice of Provenance (TaPP 2013)*. USENIX Association.
- [6] Adam M. Bates, Dave Tian, Kevin R. B. Butler, and Thomas Moyer. 2015. Trustworthy Whole-System Provenance for the Linux Kernel. In *Proceedings of the 24th USENIX Security Symposium (USENIX Security 2015)*. 319–334.
- [7] Khalid Belhajjame, Reza B'Far, James Cheney, Sam Coppens, Stephen Cresswell, Yolanda Gil, Paul Groth, Graham Klyne, Timothy Lebo, Jim McCusker, Simon Miles, James Myers, Satya Sahoo, and Curt Tilles. 2012. *PROV-DM: The PROV Data Model*. Technical Report. World Wide Web Consortium. <http://www.w3.org/TR/prov-dm/>
- [8] Sheung Chi Chan. 2019. ProvMark repository. <http://provmark2018.github.io>.
- [9] Sheung Chi Chan, Ashish Gehani, James Cheney, Ripduman Sohan, and Hassaan Irshad. 2017. Expressiveness Benchmarking for System-Level Provenance. In *Proceedings of the 8th USENIX Workshop on Theory and Practice of Provenance (TaPP 2017)*. USENIX Association.
- [10] Stephen A. Cook. 1971. The Complexity of Theorem-proving Procedures. In *Proceedings of the Third Annual ACM Symposium on Theory of Computing (STOC '71)*. ACM, New York, NY, USA, 151–158.
- [11] Martin Gebser, Benjamin Kaufmann, Roland Kaminski, Max Ostrowski, Torsten Schaub, and Marius Thomas Schneider. 2011. Potassco: The Potsdam Answer Set Solving Collection. *AI Commun.* 24, 2 (2011), 107–124.
- [12] Ashish Gehani and Dawood Tariq. 2012. SPADE: Support for Provenance Auditing in Distributed Environments. In *Proceedings of the 13th International ACM/I-FIP/USENIX Middleware Conference (Middleware 2012)*. 101–120.
- [13] Ashish Gehani, Dawood Tariq, Basim Baig, and Tanu Malik. 2011. Policy-Based Integration of Provenance Metadata. In *Proceedings of the 2011 IEEE International Symposium on Policies for Distributed Systems and Networks (POLICY 2011)*. 149–152.
- [14] Todd J. Green, Shan Shan Huang, Boon Thau Loo, and Wenchao Zhou. 2013. Datalog and Recursive Query Processing. *Foundations and Trends in Databases* 5, 2 (2013), 105–195.
- [15] Steve Grubb. 2017. Linux Audit. <http://people.redhat.com/sgrubb/audit/>.
- [16] Xueyuan Han, Thomas Pasquier, Tanvi Ranjan, Mark Goldstein, Margo Seltzer, and PDF Han. 2017. FRAPuccino: Fault-detection through Runtime Analysis of Provenance. In *Proceedings of the 9th USENIX Workshop on Hot Topics in Cloud Computing (HotCloud'17)*. USENIX Association.
- [17] Wajih Ul Hassan, Mark Lemay, Nuraini Aguse, Adam Bates, and Thomas Moyer. 2018. Towards Scalable Cluster Auditing through Grammatical Inference over Provenance Graphs. In *Proceedings of the 25th Annual Network and Distributed System Security Symposium (NDSS 2018)*.
- [18] Trung Dong Huynh, Michael O. Jewell, Amir Sezavar Keshavarz, Sanius T. Michaelides, Huanjia Yang, and Luc Moreau. 2013. *The PROV-JSON Serialization*. Technical Report. World Wide Web Consortium. <https://www.w3.org/Submission/2013/SUBM-prov-json-20130424/>
- [19] Gerwin Klein, Kevin Elphinstone, Gernot Heiser, June Andronick, David Cock, Philip Derrin, Dhammika Elkaduwe, Kai Engelhardt, Rafal Kolanski, Michael Norrish, Thomas Sewell, Harvey Tuch, and Simon Winwood. 2009. seL4: Formal Verification of an OS Kernel. In *Proceedings of the 22nd ACM Symposium on Operating Systems Principles (SOSP '09)*. ACM, New York, NY, USA, 207–220.
- [20] Luc Moreau. 2010. The Foundations for Provenance on the Web. *Foundations and Trends in Web Science* 2, 2-3 (2010), 99–241. <https://doi.org/10.1561/1800000010>
- [21] Luc Moreau et al. 2008. Special Issue: The First Provenance Challenge. *Concurrency and Computation: Practice and Experience* 20, 5 (2008), 409–418.
- [22] Kiran-Kumar Muniswamy-Reddy and David A. Holland. 2009. Causality-based Versioning. *ACM Trans. Storage* 5, 4, Article 13 (Dec. 2009), 28 pages.
- [23] Kiran-Kumar Muniswamy-Reddy, David A. Holland, Uri Braun, and Margo Seltzer. 2006. Provenance-Aware Storage Systems. In *Proceedings of the 2006 USENIX Annual Technical Conference*. 43–56.
- [24] Luke Nelson, Helgi Sigurbjarnarson, Kaiyuan Zhang, Dylan Johnson, James Bornholt, Emina Torlak, and Xi Wang. 2017. Hyperkernel: Push-Button Verification of an OS Kernel. In *Proceedings of the 26th Symposium on Operating Systems Principles (SOSP '17)*. ACM, New York, NY, USA, 252–269.
- [25] Liliana Pasquale, Dalal Alrajeh, Claudia Peersman, Thein Than Tun, Bashar Nuseibeh, and Awais Rashid. 2018. Towards forensic-ready software systems. In *Proceedings of the 40th International Conference on Software Engineering: New Ideas and Emerging Results (ICSE (NIER) 2018)*. 9–12.
- [26] Thomas Pasquier, Xueyuan Han, Mark Goldstein, Thomas Moyer, David M. Ebers, Margo Seltzer, and Jean Bacon. 2017. Practical whole-system provenance capture. In *Proceedings of the 2017 Symposium on Cloud Computing (SoCC 2017)*. 405–418.
- [27] Thomas Pasquier, Xueyuan Han, Thomas Moyer, Adam Bates, Olivier Hermant, David Ebers, Jean Bacon, and Margo Seltzer. 2018. Runtime Analysis of Whole-System Provenance. In *Conference on Computer and Communications Security (CCS'18)*. ACM, 1601–1616.
- [28] Thomas Pasquier, Matthew K Lau, Ana Trisovic, Emery R Boose, Ben Couturier, Mercè Crosas, Aaron M Ellison, Valerie Gibson, Chris R Jones, and Margo Seltzer. 2017. If these data could talk. *Scientific data* 4 (2017).
- [29] Thomas Pasquier, Jatinder Singh, David Ebers, and Jean Bacon. 2015. CamFlow: Managed data-sharing for cloud services. *IEEE Transactions on Cloud Computing* 5, 3 (2015), 472–484.
- [30] Thomas Pasquier, Jatinder Singh, Julia Powles, David Ebers, Margo Seltzer, and Jean Bacon. 2018. Data provenance to audit compliance with privacy policy in the Internet of Things. *Personal and Ubiquitous Computing* 22, 2 (2018), 333–344.
- [31] Devin J. Pohly, Stephen E. McLaughlin, Patrick D. McDaniel, and Kevin R. B. Butler. 2012. Hi-Fi: collecting high-fidelity whole-system provenance. In *Proceedings of the 28th Annual Computer Security Applications Conference (ACSAC 2012)*. 259–268.
- [32] Joerg Thalheim, Pramod Bhatotia, and Christof Fetzer. 2016. Inspector: Data Provenance using Intel Processor Trace (PT). In *Proceedings of IEEE International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 25–34.
- [33] Wei Wang and Thomas E. Daniels. 2008. A Graph Based Approach Toward Network Forensics Analysis. *ACM Trans. Inf. Syst. Secur.* 12, 1, Article 4 (Oct. 2008), 33 pages.
- [34] Jieming Zhu, Pinjia He, Qiang Fu, Hongyu Zhang, Michael R. Lyu, and Dongmei Zhang. 2015. Learning to Log: Helping Developers Make Informed Logging Decisions. In *Proceedings of the 37th IEEE/ACM International Conference on Software Engineering (ICSE 2015)*. 415–425.