

# Efficient Differencing of System-level Provenance Graphs

Yuta Nakamura  
DePaul University  
Chicago, IL, USA  
ynakamu1@depaul.edu

Iyad Kanj  
DePaul University  
Chicago, IL, USA  
ikanj@depaul.edu

Tanu Malik  
DePaul University  
Chicago, IL, USA  
tanu.malik@depaul.edu

## ABSTRACT

Data provenance, when audited at the operating system level, generates a large volume of low-level events. Current provenance systems infer causal flow from these event traces, but do not infer application structure, such as loops and branches. The absence of these inferred structures decreases accuracy when comparing two event traces, leading to low-quality answers from a provenance system. In this paper, we infer nested natural and unnatural loop structures over a collection of provenance event traces. We describe an ‘unrolling method’ that uses the inferred nested loop structure to systematically mark loop iterations. Our loop-based unrolling improves the accuracy of trace comparison by 20-70% over trace comparisons that do not rely on inferred structures.

## CCS CONCEPTS

• **Information systems** → **Data management systems**; • **Data management** → *Data structures*; • **Data structures** → Data access methods; • **Data access methods** → Unidimensional range search.

## KEYWORDS

OS-level auditing, data provenance querying, graph and sequence alignment, loop identification, unnatural loops

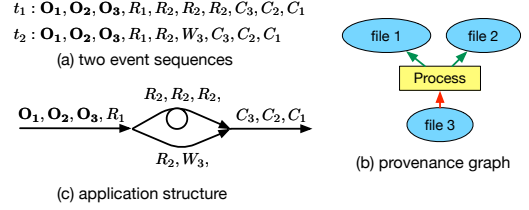
### ACM Reference Format:

Yuta Nakamura, Iyad Kanj, and Tanu Malik. 2023. Efficient Differencing of System-level Provenance Graphs. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management (CIKM '23)*, October 21–25, 2023, Birmingham, United Kingdom. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/3583780.3615171>

## 1 INTRODUCTION

Data provenance is a record of the origin and evolution of data in a program or application. Audited provenance is vital for establishing reproducibility [4, 7], determining security breaches [14], and for diagnostics [6]. In its raw form, provenance data is simply a series of events, such as system calls or function entries and exits, generated during an application execution. Provenance systems [9, 18, 19] create a lineage graph using some of the events for lineage queries.

Consider a simple sequence of events (labeled  $t_1$  in Figure 1(a)) that involves the system calls `open`, `read`, and `close`. The nodes and edges of the provenance graph (in Figure 1(b)) are obtained from



**Figure 1: Low-level event traces, used to generate a provenance graph, omit out some events that represent application structure.**

some events (in bold), such as the opening of a file by a process for reading or writing, as they determine the causal data flow between the process and the file. Other events, such as reading the file or closing the file, do not introduce new edges into the provenance graph.

Current systems [9, 18, 19] simply link other events to one of two nodes in the causal edge graph. These event sequences, however, hold vital information about the program or application structure. For example, if the same application is run twice with different input parameters, and each time it opens the same file for reading but in one trace ( $t_1$  in Figure 1(a)) simply reads the file, and in the other ( $t_2$  in Figure 1(a)) also writes to the file, it exhibits a change in application execution behavior due to a different control flow path being exercised on the changed input. In the same way, the repetition of a read call from a program location indicates a loop and branch application structure (Figure 1(c)). The determination of this application structure is necessary for improving the fidelity of the provenance graph or for comparing two application executions.

In this paper, we infer the application structure using event traces generated by multiple executions of an application. We infer two types of application structures: loops and branches. Loops have strongly connected events, while branches have mandatory events followed by optional events. In general, the repeated occurrence of uniquely identifiable events indicates the occurrence of a loop. However, loops may be nested, may have single or multiple points of entry, or may have branches. Thus, we not only need to infer the existence of a loop, but also to determine the type of loop structure to improve the fidelity of the provenance graph.

We use the inferred loop structure to compare any two application event traces. Our research indicates that when comparing two traces, identifying differences that align with loop boundaries leads to more precise reporting of differences compared to a basic edit-distance comparison of traces. In Section 3, we describe an ‘unrolling’ based method to obtain loop boundary locations where traces may have diverged or converged, and show the accuracy of our approach in Section 4.

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 republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
CIKM '23, October 21–25, 2023, Birmingham, United Kingdom  
© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 979-8-4007-0124-5/23/10...\$15.00  
<https://doi.org/10.1145/3583780.3615171>

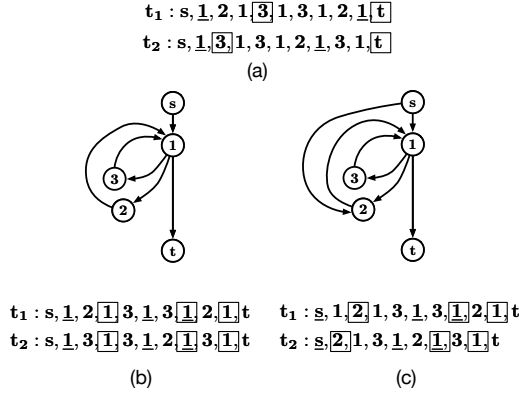


Figure 2: (a) traces compared via edit-distance (b) traces compared per natural loop graph (c) traces compared per unnatural loop (entry from both 1 and 2) graph.

## 2 MOTIVATING EXAMPLE

We consider an example to show how inferred loops improve accuracy when comparing two traces. Let us assume a user executes an application  $\mathcal{A}$  two times resulting in two event traces  $t_1$  and  $t_2$ . Each trace generates a sequence of event identifiers that have been disambiguated and made uniquely identifiable through respective event properties<sup>1</sup>.

Figure 2(a) shows two such event traces in which some events repeat. If we compare the linear traces, using edit distance and without accounting for the underlying graph structure, we highlight differences starting at  $\underline{\phantom{x}}$ , i.e., when prefixes do not match, and ending at  $\boxed{\phantom{x}}$ , i.e., when a subsequent same event is witnessed again in the other trace (in general, order of comparison matters). As the Figure 2(a) shows, these event locations of  $\underline{\phantom{x}}$  and  $\boxed{\phantom{x}}$  are determined without knowledge of the graph application structure of traces. If, however, the traces are represented as graphs and the graph representation of Figure 2(b) is accounted for, event 3 is not the location where differences end. As the graph shows, no two different paths merge at event 3. Indeed, the two traces both diverge and converge at event 1. Due to this change in location from 3 to 1, other distinct differences are found that are aligned according to the graph structure. In total, a higher number of differences—2  $\boxed{\phantom{x}}$ s in Figure 2(b) instead of 3—are highlighted than when considering the edit distance.

The example only assumes a simple nested loop structure. In general, provenance traces may result from unnatural loops (more than one entry point into the loop) as demonstrated in Figure 2(c). Determining the loop hierarchy, i.e., which loop is nested within which loop, and being consistent about where loops start and end, i.e., if the loop starts from event 1 or 2, is important to correctly compare event traces.

## 3 INFERRING APPLICATION STRUCTURE

We use notions of loop and loop tree as defined for control-flow analysis [3, 11]. We add provenance trace meta information to loop tree and use it to locate divergent/convergent points by “unrolling” the graph. We note a branch structure is a simplified loop and, for simplicity, restrict to general loops.

<sup>1</sup>We omit the disambiguation details as they are beyond the scope of the current paper. More details are in [16].

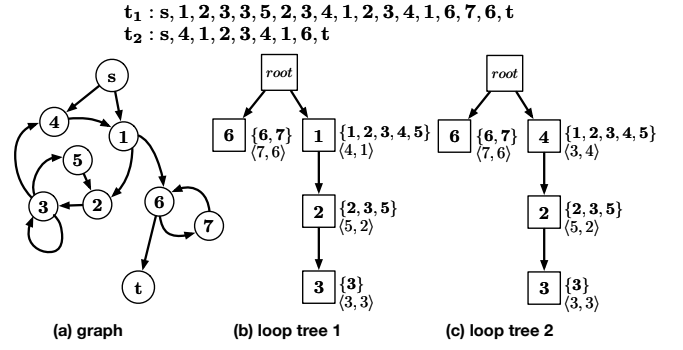


Figure 3: The graph obtained from two traces  $t_1$  and  $t_2$ , and its two possible loop trees. In loop trees, the node vertices (in squares) are entry points, numbers in the curly brackets correspond to graph vertices, and edges in  $\langle \phantom{x} \rangle$  are back edges.

### 3.1 Inferring Loop Structures

When traces are given, we can create a graph out of the traces. Let  $t = (v_1, \dots, v_s)$  be a trace and  $G_t = (V_t, E_t)$  its associated graph, where  $v_i \in V_t, \forall i \in \{1, \dots, s\}$ , and  $(v_i, v_{i+1}) \in E_t, \forall i \in \{1, \dots, s-1\}$ .

**Definition 3.1.** Let  $Q$  be a subgraph of a graph  $P$ . A vertex  $v \in V(Q)$  is an *entry point* for  $Q$  if  $v$  has at least one incoming edge that is not in  $Q$  (i.e., is in  $E(P) - E(Q)$ ).

**Definition 3.2.** We define the *loops* of  $G$  inductively as follows.

- (1) For a graph  $Q$  such that  $V(Q) = \{v\}$ ,  $Q$  is a loop iff  $(v, v)$  is an edge in  $Q$  (i.e.,  $(v, v)$  is a self-loop).
- (2) Suppose, inductively, that the loops have been defined for any graph of vertex-cardinality less than  $|V(G)|$ . If  $G$  is strongly connected, then  $G$  is a loop. Moreover, for every entry vertex  $v \in V(G)$ , the loops of  $V(G) - \{v\}$  are loops of  $G$ .

Thus, one can identify the type of loop by considering its entry points and the strongly connected component (SCC) it belongs to. If a loop has a single entry point, it is a *natural* loop, and if it has more than one entry point, then it is an *unnatural* loop. Strongly connected components are determined as per [21].

**Definition 3.3.** A *loop tree*,  $\mathcal{T}$ , of a graph  $G$  is defined as follows. The root of  $\mathcal{T}$  is a dummy node whose children are nodes representing the SCC's of  $G$ . The children of a node in  $\mathcal{T}$  are defined recursively as follows. For each node  $\alpha$  in the tree, arbitrarily fix an entry point  $v_\alpha$  of the subgraph  $G_\alpha$  (of  $G$ ) representing node  $\alpha$ . (Note that if  $G$  itself is a SCC, then it has an entry point.) For each SCC  $C$  of  $G_\alpha - v_\alpha$ , create a child of  $\alpha$  that represents  $C$ .

Figure 3 depicts the loop tree. Let  $G$  be the graph created from  $T_1$  and  $T_2$ . We perform a depth-first search (DFS) on  $G$ ; recall that a DFS results in a DFS forest, and that the back edges of  $G$  (with respect to the DFS forest) are the edges that go between descendants and ancestors in the forest. The SCCs of  $G$  are  $\{1, 2, 3, 4, 5\}$  and  $\{6, 7\}$ . The entry points of the SCC  $\{1, 2, 3, 4, 5\}$  can be 1 or 4, and we pick 1 for the loop tree in Figure 3(b). For  $\{6, 7\}$ , we have to choose 6. Then, in the subgraph  $\{1, 2, 3, 4, 5\} - \{1\}$ , one can find the SCC  $\{2, 3, 5\}$  and 2 as an entry point. Finally, from  $\{2, 3, 5\} - \{2\}$ , we will

get SCC  $\{3\}$  and its entry point 3. If we had chosen 4 as the entry point, then we get an alternate loop tree with a different back edge (Figure 3(c)).

### 3.2 Creating an Unrolled Graph for Trace Comparison

To compare two traces, we first define the points of divergence and convergence.

Let  $T_1, T_2$  be two traces, and let  $G_{T_1}, G_{T_2}$  be the graphs induced by  $T_1$  and  $T_2$ , respectively. Also, for a trace  $T = (v_1, \dots, v_k)$  in  $G$  and for a vertex  $v_i \in T$ , define  $prev(v_i) = v_{i-1}$  if  $i > 1$  and  $prev(v_i) = \perp$  if  $i = 1$  (where  $\perp$  denotes nil, i.e., indicating that  $v_1$  does not have a previous vertex in  $T$ ). Similarly, for  $v_i \in T$ , define  $next(v_i) = v_{i+1}$  if  $i < k$ , and  $\perp$  otherwise (i.e.,  $i = k$ ).

**Definition 3.4 (Points of Divergence and Convergence).** Let  $T_1 = (v_1, \dots, v_r)$  and  $T_2 = (w_1, \dots, w_s)$  be two traces. A vertex  $v_i \in T_1$ ,  $i \geq 1$ , is said to be a *point of divergence* for  $T_1$  and  $T_2$  if there exists  $w_j \in T_2$ ,  $j \geq 1$ , such that  $v_i = w_j$  and  $next(v_i) \neq next(w_j)$ . A vertex  $v_i \in T_1$ ,  $i \geq 1$ , is said to be a *point of convergence* for  $T_1$  and  $T_2$  if there exists  $w_j \in T_2$ ,  $j \geq 1$ , such that  $v_i = w_j$  and  $prev(v_i) \neq prev(w_j)$ .

**3.2.1 The case where  $G = G_{T_1} \cup G_{T_2}$  is a DAG.** If  $G = G_{T_1} \cup G_{T_2}$  is a directed acyclic graph, then  $next(v_i)$ ,  $next(w_j)$  and  $prev(v_i)$ ,  $prev(w_j)$ , on each respective trace are uniquely defined.

**3.2.2 The case where  $G = G_{T_1} \cup G_{T_2}$  has loops.** In the case of loops, the  $next$  and  $prev$  are not uniquely defined and we need a loop tree to determine where a loop starts and ends. Let  $\mathcal{T}$  be a loop tree of  $G$  and let  $F_G$  be a depth-first search (DFS) forest resulting from performing DFS on  $G$ . We define an extended graph,  $\Gamma_G$ , from  $G$  by repeatedly “unrolling” every loop corresponding to a node in  $\mathcal{T}$  as explained next.

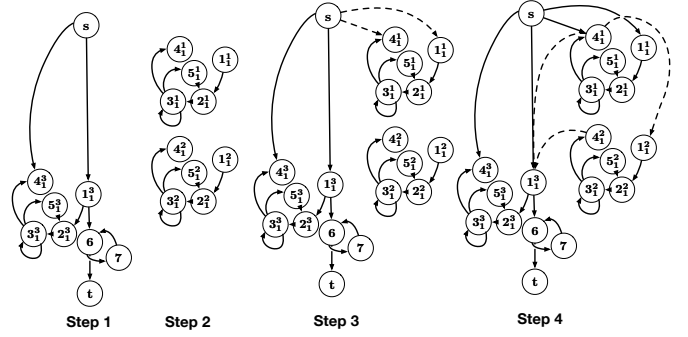
- (i) Initialize  $\Gamma_G$  to  $G$ . Let  $\alpha$  be a node in  $\mathcal{T}$  of entry point  $v_\alpha$ , and let  $G_\alpha$  be the subgraph of  $G$  corresponding to  $\alpha$ . Let  $e_1, \dots, e_k$ , where  $e_i = (u_i, v_\alpha)$ , for  $i = 1, \dots, k$ , be the back edges with multiplicities (i.e., possibly with repetition), with respect to  $F_G$ , incoming to  $v_\alpha$  in  $G_\alpha$ , in the order in which they are traversed by  $T_1$  and then by  $T_2$ .
- (ii) Let  $G_\alpha^-$  be  $G_\alpha$  with all these back edges removed. We remove the back edges from  $\Gamma_G$  (step 1), and add  $k$  copies of  $G_\alpha^-$  to  $\Gamma_G$ . (We will denote them  $G_\alpha^{-1}, G_\alpha^{-2}, \dots, G_\alpha^{-k}$ ) (step 2). Create edges from the sources of the entry edges in  $G$  to all of the entry points in  $G_\alpha^{-1}$  (not only  $v_\alpha^1$ ) (step 3). Let  $v_\alpha^{k+1} = v_\alpha$  and let  $u_\alpha^i, v_\alpha^i$  denote the copies of vertices  $u_i$  and  $v_\alpha$ , respectively, for  $i = 1, \dots, k$ , in the  $i$ -th copy of  $G_\alpha^-$ .
- (iii) For each back edge  $e_i = (u_i, v_\alpha)$ , for  $i = 1, \dots, k$ , we add an edge in  $\Gamma_G$  from  $u_\alpha^i$  in  $\Gamma_G$  to  $v_\alpha^{i+1}$  in  $\Gamma_G$ . We also add edges from  $u_\alpha^i$  to  $v_\alpha^{k+1}$  for  $i = 1, \dots, k$  (step 4).

Consider  $\Gamma_G$  after unrolling all the loops in  $\mathcal{T}$ . We have:

**THEOREM 3.1.**  $\Gamma_G$  is a DAG.

The proof (included in our technical report [15]) follows by showing that no cycle in  $G$  remains in  $\Gamma_G$ , since our creation of  $\Gamma_G$  does not introduce a cycle that is not a cycle in  $G$ .

Since  $\Gamma_G$  is a DAG, we can now use the same definition of points of divergence/convergence but on  $\Gamma_G$ .



**Figure 4: Unrolling of one loop in Figure 3 (a).** The back edge is  $\langle 4, 1 \rangle$ , so at first the back edge of  $G$  are omitted and  $k$  copies of  $G_\alpha^-$  are created (Step 1, 2). For the 1st iteration of the loop, all the entry edges from original  $G$  to entry nodes in  $G_\alpha^-$  are created. (Step 3). For all the other iterations, create an edge from  $u_\alpha^{i-1}$  to  $v_\alpha^i$  and to  $v_\alpha^{k+1}$  for  $\forall i \in \{2, \dots, k\}$  (Step 4).

## 4 EXPERIMENTAL EVALUATION

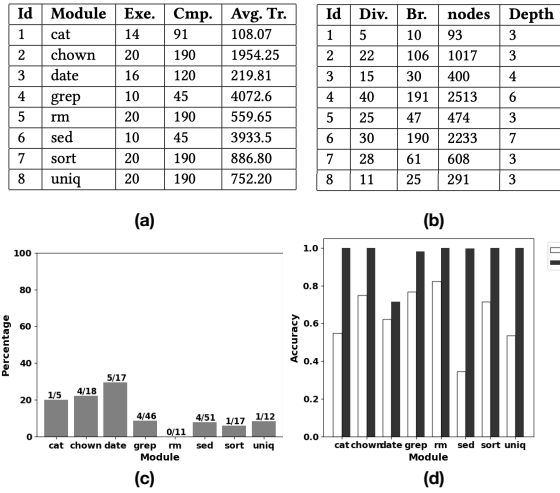
**Provenance Collection.** We used the the Pinsystem [10, 20] to audit provenance traces. Pin uses a combination of system calls and function calls to audit application execution. We consider eight different applications as described in table of Figure 5(a). These applications are taken from Coreutils [2] and SIR [1]. The objective was to experiment with realistic, widely-used programs that have simple source code so, if necessary, we can validate the actual application structure. Raw provenance is collected by executing the application with different input parameters. Tables in Figure 5(a) and Figure 5(b) describes properties of collected traces, the obtained graph, and loop tree from the traces.

### 4.1 Inferred Loops

To measure fidelity, we report, in Figure 5(c), the number of inferred natural and unnatural loops for each module. The total number of loops are those loops determined over the subgraph induced by the traces reported in column #3 of the table in Figure 5(a). In general, the application may have more loops. As Figure 5(c) shows, depending on the module, the percentage of unnatural loops can be as high as 30%, showing the necessity to determine both types of loops.

### 4.2 Precision of Trace Comparison

To compute the precision of a difference-based lineage query, we compare two provenance traces, one in which application structure is not inferred, and one in which it is inferred. We use an edit distance measure for comparison. To report the results, we compute if the divergence points were the same or different across the two methods. We assign a measure of 0 if points of divergence were totally different, and 1 if they are the same. Thus higher the number, the more similar the points of divergence are. As Figure 5(d) shows, several points of divergence (obtained by comparing two traces with unnatural loops) are ignored if the application structure is not considered (white bar), and a few points are ignored if unnatural loops are not considered (black bar). Comparing the



**Figure 5: (a) Trace Details (Exe. = # of executions, Comp. = # of comparisons, Avg. Tr = Average Trace Length) (b) Trace-induced graph and loop tree details (Div. = total # of divergences across all traces, Br. = total # of branches across all traces, Nodes = total # of nodes across all traces, Depth = total depth of loop tree.) (c) # of Loop structures (d) Precision of Trace Comparison.**

results obtained in Figure 5(d) with the actual ground truth requires a finer granularity of the trace and employing dynamic tainting procedures. Comparing with dynamic tainting approaches such as DynamicRio [5] is part of our future work.

## 5 RELATED WORK

Several systems [9, 18, 19] collect provenance by monitoring executions at the operating system. Some edit-based methods [4, 16, 22] to compare provenance traces have recently been developed. These methods either do not apply to system call traces. In [8], the provided tools and utilities for comparing provenance are based on edit-distance [8] and do not account for application structure.

We leverage loop identification methods from control flow analysis. Many papers have looked into nested loops, but often they lack in the consideration of unnatural loops [12, 13]. [23] presents a loop identification algorithm on control-flow graphs for both natural and unnatural loops. We identify loops on the restricted control flow graph that is witnessed by lineage traces and show that we must choose from a set of loop trees to unroll the graph. While loop identification has been used to compress traces [17], ours is the first method to show the use of the inferred loop structures to unroll the graph to compute accurate differences.

## 6 CONCLUSION

In this paper, we improved fidelity of system-level provenance graphs by considering application structure. We showed that inferring loops in vital for accurately determining differences across traces and provided a method to unroll a graph based on types of loops so divergence and convergence points are accurately obtained.

Our results show comparison accuracy upto 70% over simple edit distance measures. In the future, we plan to identify loop identification in traces from parallel applications to determine how differing inner structure of parallel programs affects performance.

## ACKNOWLEDGEMENTS

This work is supported by the National Science Foundation under grants CNS-1846418 and by the National Aeronautics Space Agency under grant AIST-21-0095-80NSSC22K1485.

## REFERENCES

- [1] 2006. Software-artifact Infrastructure Repository. <https://sir.csc.ncsu.edu>.
- [2] 2016. Coreutils - GNU core utilities. <http://www.gnu.org/software/coreutils/>.
- [3] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. 2006. *Compilers: Principles, Techniques, and Tools* (2 ed.). Pearson Education, Inc.
- [4] Zhuwei Bao, Sarah Cohen-Boulakia, Susan B Davidson, Anat Eyal, and Sanjeev Khanna. 2009. Differencing provenance in scientific workflows. In *ICDE*.
- [5] Derek Bruening and Saman Amarasinghe. 2004. Efficient, transparent, and comprehensive runtime code manipulation. (2004).
- [6] Ang Chen, Yang Wu, Andreas Haeberlen, Wencho Zhou, and Boon Thau Loo. 2016. The good, the bad, and the differences: Better network diagnostics with differential provenance. In *SIGCOMM*.
- [7] Juliana Freire, Cláudio T. Silva, Steven P. Callahan, Emanuele Santos, Carlos E. Scheidegger, and Huy T. Vo. 2006. Managing Rapidly-Evolving Scientific Workflows. In *IPAW*.
- [8] Ashish Gehani, Raza Ahmad, Hassaan Irshad, Jianqiao Zhu, and Jignesh Patel. 2021. Digging into big provenance (with SPADE). *Commun. ACM* 64, 12 (2021), 48–56.
- [9] Ashish Gehani and Dawood Tariq. 2012. SPADE: Support for Provenance Auditing in Distributed Environments. In *Middleware*.
- [10] Intel. 2012. Pin: A Dynamic Binary Instrumentation Tool. <https://software.intel.com/content/www/us/en/develop/articles/pin-a-dynamic-binary-instrumentation-tool.html>.
- [11] Ken Kennedy and John R. Allen. 2001. *Optimizing Compilers for Modern Architectures: A Dependence-based Approach*. Morgan Kaufmann Publishers Inc.
- [12] Alain Ketterlin and Philippe Clauss. 2008. Prediction and Trace Compression of Data Access Addresses through Nested Loop Recognition. In *Proceedings of the 6th Annual IEEE/ACM International Symposium on Code Generation and Optimization*. 94–103.
- [13] Kobayashi. 1984. Dynamic Characteristics of Loops. *IEEE Trans. Comput.* C-33, 2 (1984), 125–132. <https://doi.org/10.1109/TC.1984.1676404>
- [14] Kyu Hyung Lee, Xiangyu Zhang, and Dongyan Xu. 2013. High Accuracy Attack Provenance via Binary-based Execution Partition. In *NDSS*. 16.
- [15] Kanj Iyad Nakamura, Yuta and Tanu Malik. 2023. *Efficient Differencing of System-level Provenance Graphs*. Technical Report. School of Computing, DePaul University. <https://dice.cs.depaul.edu/pdfs/pubs/C33.pdf>
- [16] Yuta Nakamura, Tanu Malik, and Ashish Gehani. 2020. Efficient Provenance Alignment in Reproduced Executions. In *12th International Workshop on Theory and Practice of Provenance (TaPP 2020)*.
- [17] Michael Noeth, Prasun Ratn, Frank Mueller, Martin Schulz, and Bronis R De Supinski. 2009. Scalatrace: Scalable compression and replay of communication traces for high-performance computing. *J. Parallel and Distrib. Comput.* 69, 8 (2009), 696–710.
- [18] Thomas Pasquier, Xueyan Han, Mark Goldstein, Thomas Moyer, David Eysers, Margo Seltzer, and Jean Bacon. 2017. Practical whole-system provenance capture. In *SoCC*.
- [19] Devin J Pohly, Stephen McLaughlin, Patrick McDaniel, and Kevin Butler. 2012. Hi-Fi: collecting high-fidelity whole-system provenance. In *Proceedings of the 28th Annual Computer Security Applications Conference*. 259–268.
- [20] Vijay Janapa Reddi, Alex Settle, and et. al. 2004. Pin: a binary instrumentation tool for computer architecture research and education. In *Workshop on Computer architecture education*.
- [21] Robert Endre Tarjan. 1972. Depth-first search and linear graph algorithms. *SIAM J. Comput.* 1, 2 (1972), 146–160.
- [22] P. Thavasimani, J. Cala, and P. Missier. 2019. Why-Diff: Exploiting Provenance to Understand Outcome Differences From Non-Identical Reproduced Workflows. *IEEE Access* (2019).
- [23] Tao Wei, Jian Mao, Wei Zou, and Yu Chen. 2007. A new algorithm for identifying loops in decompilation. In *Static Analysis: 14th International Symposium, SAS 2007, Kongens Lyngby, Denmark, August 22–24, 2007. Proceedings 14*. Springer, 170–183.