D-031

# Preserving Integrity and Confidentiality of a Directed Acyclic Graph Model of Provenance

Amril Syalim†        Takashi Nishide†        Kouichi Sakurai†

## 1. Introduction

In a system where we need to understand the processes that have been executed to produce a result we need to record provenance of the execution. By recording provenance we may trace who have contributed to the creation of the result. This feature is very important whenever we need to verify the process of result's creation, for example in a distributed system (i.e. a grid system), where a result may be produced by many parties in different computers. Another real life example is in a hospital, a medicine prescription may be created by a doctor by collaborating with other doctors.

For a sequential execution of processes, provenance can be represented in a form of chain [1,2]. A more expressive model that is suitable for a parallel execution is a directed graph model where nodes in the graph represent processes and the edges represent relationships between the processes (nodes) [3,4]. Because provenance is tightly associated with time, many models of provenance take the form of a directed acyclic graph (DAG) [5].

In this paper, we are focusing on securing a directed acyclic graph model of provenance in terms of integrity and confidentiality. To ensure integrity of the provenance graph (i.e. nodes and edges) we need to assure immutability and non-repudiation properties of each node and edge in the provenance graph. The contributors (i.e. the people or processes that contribute in the provenance graph) can not cheat for any purposes. The other parties in the provenance system (i.e. the manager of the provenance graph that we refer in this paper as the provenance owner), although powerful enough to manage access to provenance, also can not cheat (i.e by changing the provenance graph) without being detected.

We propose a method to protect integrity of provenance by employing digital signature. Using this method, the contributors and provenance owner both sign the provenance's nodes and edges. To alter the nodes and edges without detected needs collusion from the two parties which means to repeat the process execution from beginning. To support confidentiality of the provenance graph we need to define access control model to provenance and how to enforce the access control model. Provenance should only be accessed by the person who has the right to access, for example an auditor who need to audit the process. The system should support restricting access to only some parts of the provenance.

Many access control models employ grouping mechanism to improve efficiency and security (i.e. by using groups, roles, security levels/compartments). We propose a grouping mechanism for access control to provenance by utilizing two

†Department of Informatics, Kyushu University

grouping methods: grouping of entities in the provenance graph based on paths and grouping entities based on compartments. Grouping by paths is useful because the auditors who audit the process should be interested in the causal relationship in the provenance graph.

However access control by paths alone is not expressive to enforce more specific policies (i.e. an auditor only can access a part of nodes/edges in the paths). We complement the paths-based access control with a compartment-based access control so that we can enforce such policies. By using a compartment-based access control, each node is assigned with a compartment and the provenance owner grants access to the nodes in a compartment by granting access to that compartment.

## 2. Integrity Mechanism: Digitally Signing the Provenance Graph

An example of provenance graph with six contributors is shown in the Figure 1. The Figure 1 shows that to produce the final result, the contributor C5 uses the outputs of contributors C1 and C2 while contributor C6 uses the output of contributors C3 and C4. Contributor C7 uses the output of C5 and C6 which later used by C8 and C9. The final process is executed by C10 that processes the outputs of C8 and C9. After each process is executed and the provenance of the process (i.e. node) is created/generated, the provenance is stored in the provenance database. The other papers call the provenance database as a provenance store.
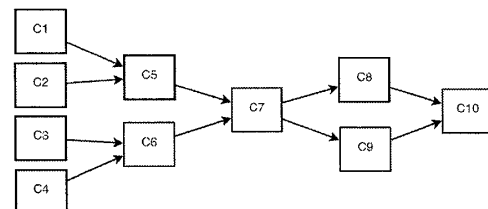


Figure 1. Provenance Graph

We identify three groups of active entities involved in a provenance system: provenance owners, contributors, and auditors. A provenance owner is the owner of provenance that mediates the provenance recording process and manages access to the provenance. The contributors are the people who execute process and contribute the results. Auditors are the people who need to access the provenance graph, for example for reviewing or auditing the process's execution.

Provenance is recorded after each process is executed by the contributor. In a distributed system, before executing the distributed process, a worflow (i.e. a distributed execution plan) should be defined and sent to the provenance owner. The process to create a workflow may involve some or all of contributors. Based on the workflow, provenance owner sends

each contributor information that is needed by the contributor to execute each process in the workflow (i.e inputs of the process). After a contributor execute a process, the contributor should produce outputs which we refer as a document. The provenance of the document is documentation of process execution to produce the document. The provenance can be automatically generated by the system where the contributor execute the process or manually created by the contributor.

After execution of a process, the document and provenance of the document are sent to the provenance owner which later record them as a node in the provenance database. The provenance owner may also send the document to contributors that need the documents for their inputs.

After the provenance is recorded, there are some possible integrity problems with provenance. We identify four main problems: repudiation, alteration, deletion, and addition. A contributor may deny that she/he has contributed the document and its provenance. The document and provenance (i.e. nodes) may be altered by an attacker so that they do not reflect original process. Attacker may also delete a node or add a fake node.

The basic idea of the digital signature mechanism is whenever a provenance of a document is recorded, both of the contributor and provenance owner sign the document and the provenance before storing the provenance to the database. Whenever a contributor uses an output document of other contributor as an input, the contributor should create the hash/checksum of the input and store them as a provenance of the process executed by the contributor.

We assume that each contributor, auditor and provenance owner has a pair of public key and private key and each party can retrieve the public keys of the other parties securely. The private keys can only be accessed by the owner of the key. Let $D_n$ is the document created by a contributor identified by n and $P_n$ is a provenance of the document. The function $H(D_n)$ is a function that produce hash value of $D_n$. The function $S_n$ is a signing function where $S_n(D_n)$ is a function that produce digital signature of contributor n to document $D_n$. N is the number of inputs used by a process to produce a document.

If a contributor n needs to use a document (i.e. $D_{n-1}$) produced by another contributor (i.e. contributor n-1) as input, before the contributor n executes the process, the provenance owner sends the input that has been signed by the provenance owner and the another contributor: $S_o(S_{n-1}(D_{n-1}))$. After verifying the document and the signatures, the contributor n execute the process. The contributor n signs the result $D_n$, its provenance $P_n$ and hash of the input $H(D_{n-1})$. The signed result, its provenance and hash of the input is $S_n(P_n, D_n, H(D_{n-1}))$. The contributor sends them to the provenance owner. The provenance owner signs them and stores them in the database.

## 3. Confidentiality Mechanism: Path-based Access Control and Encrypting the Provenance Graph

To protect confidentiality of provenance we need to prevent confidential provenance information be accessed by unauthorized people accessing the system. However, the system should also support authorized access to provenance (i.e. authorized auditors

who need to access provenance to do audit and verify the process of object creation). We propose an access control model based on path on the provenance graph. The arguments of our proposal is that an auditor normally needs to access all nodes that have a path to the result because the nodes have causal relationship to the result. We believe that this model is more efficient and comfortable because the provenance owner can easily create access based on paths in the provenance graph.

However, by using path-based only access control, we can not create a more expressive policy (for example an auditor can only access a part of the paths). We combine path-based policy with another access policy based on compartments. Compartments define separation between nodes in different security level/classes and the auditors that can access those compartments.

We propose to implement the access control model by using cryptographic mechanisms (i.e. encryption). This method is especially important if we store the provenance in an untrusted server (i.e. the provenance owner wants to outsource the storage of provenance to a third party who may be not trusted). This method can also be used if the provenance owner wants to implement cryptographic-based access control (where the data is encrypted and access rights are granted by giving the encryption keys). The idea of our implementation for path-based access control is to encrypt the nodes and store the encryption keys in the children of the nodes.

Let $P_n$ is the node that has been signed by the contributor n and the provenance owner o and let $E_k(P_n)$ is an encryption function that encrypt $P_n$ with private key k. To encrypt the node $P_n$, the provenance owner define compartment of the node and find the parent nodes. The provenance owner retrieves the key associated with the compartment $K_C$, the keys to encrypt the parent nodes $K_{n-1}$ and the key to encrypt the grandparent node $K_{k-1}$. The provenance owner generates two random keys: node's key $K_n$ and parent-key's key $K_k$ and store the keys in a key database managed by the provenance owner. The provenance owner encrypts the node $P_n$ with key $K_C$. Then the provenance owner re-encrypts the node with the key $K_n$. After that the provenance owner encrypts the keys $K_{n-1}$ and $K_{k-1}$ with parent-key $K_k$. Encrypted form of the node is $E_{K_n}(E_{K_C}(P_n))||E_{K_k}(K_{n-1}||K_{k-1})$. The provenance owner stores encrypted form of the node in the provenance database.

## References

[1] Hasan, R., Sion, R., Winslett, M.: Preventing history forgery with secure provenance. ACM Transactions on Storage 5(4), 12:1–12:43 (2009)

[2] Hasan, R., Sion, R., Winslett, M.: The case of the fake picasso: Preventing history forgery with secure provenance. FAST 2008.

[3] Moreau, L., Freire, J., Futrelle, J., McGrath, R.E., Myers, J., Paulson, P.: The open provenance model: An overview. IPAW 2008.

[4] Bowers, S., McPhillips, T., Lud"ascher, B., Cohen, S., Davidson, S.B.: A model for user-oriented data provenance in pipelined scientific workflows. IPAW 2006.

[5] Braun, U., Shinnar, A., Seltzer, M.I.: Securing provenance. HotSec (2008)