# Towards Eidetic Blockchain Systems with Enhanced Provenance

Shlomi Linoy
*Faculty of Computer Science*
*University of New Brunswick*
*Fredericton, Canada*
*Email: slinoy@unb.ca*

Suprio Ray
*Faculty of Computer Science*
*University of New Brunswick*
*Fredericton, Canada*
*Email: sray@unb.ca*

Natalia Stakhanova
*Department of Computer Science*
*University of Saskatchewan*
*Saskatoon, Canada*
*Email: natalia@cs.usask.ca*

*Abstract*—Modern blockchain systems with smart contract support are continuing to be rapidly adopted across various industry sectors and are increasingly used to manage valuable assets. As the size and complexity of smart contract applications increases, so are the coding errors, exploit potential, and regulation requirements. For these reasons, it has become necessary to efficiently manage the system's historic execution information, or *provenance*, to enable efficient analysis. Existing approaches facilitate historic data access, however, they do not support tracking what initiated the changes or why the states mutated. To address this, we propose a system that enables efficient management of historic smart contracts calls, their parameters, and the blockchain state before and after a call. We further explore how querying this historic data in different granularity levels can facilitate the analysis of a use case example comprised of multiple smart contract calls across different entities.

*Index Terms*—Blockchain, provenance, call graph, defect analysis, debugging, regulatory requirements, audit

## I. INTRODUCTION

Blockchain systems increasingly capture the attention of academia as well as industry at a rapid rate. The adoption of blockchain technologies is expanding in various industry sectors such as healthcare [3], IoT [4], and securities trading [5]. With a projected market size of $4.19 billion in 2020 and $162.84 billion in 2027 [13]. Modern blockchains such as Ethereum [1] and Hyperledger [2] make use of smart contracts (abbreviated as 'contracts'), which are programs that operate on the blockchain's current global state and produce a new global state. Each contract can be executed by a client, a contract, or can periodically run using a scheduler. As the technology mature, companies are increasingly using blockchain applications to manage valuable assets, including cryptocurrencies, securities, real estate and valuable tangible assets. Hence, it is important that contracts are free of defects. Issues such as coding errors, malicious code and non-compliance with regulations can result in impactful financial repercussions. For instance, it has been reported [14] that a software bug involving the replacement of += operation with =+ lead to the loss of assets worth $800,000. Another incident [15] involved an attacker exploiting a defect in contract code, which resulted in a $80 million loss. A blockchain system that can efficiently manage the historic information, or *provenance*, of both data and contract execution flow can be used to facilitate forensic analysis of contracts' malicious activities, defects, or support auditing. However, existing blockchain-based provenance so-lutions, such as [8], do not track why or how the contract states evolved, for instance, which contract executions mutated these states. Therefore, they are not suitable for root-cause analysis of blockchain transactions' changes.

Understanding the change history of data has been extensively studied in the database systems community and is referred to as *data provenance*. In data provenance each row in the output of a single query (which is possibly comprised of sub queries) is annotated with the input tuples that derived it. Cheney et al. [11] introduced three types of data provenance for a specific output tuple in a query result: *Why-provenance* - the set of minimal input tuples that contributed to the output tuple; *How-provenance* - specifies how the output tuple was generated from the minimal input tuples; and *Where-provenance* - maps the specific output tuple's fields to the input tuples' fields. Glavic et al. [6] presented a system that supports all three types of provenance by using query rewrites to annotate the output tuples with the corresponding provenance. While data provenance in general is mostly concerned with content, *workflow provenance* looks at the flow of execution and as such is derived from multiple components, each of which has its own configuration parameters, which receives, processes, and forwards data from/to other components. Miao et al. [12] propose a system that collects workflow provenance in a collaborative workflows' environment.

Due to its potential, researchers explored using blockchain as a tamper-proof, fault-tolerant, distributive, and decentralized database for storing provenance information [16] [17]. Interest in blockchain-based data provenance has been steadily growing. The blockchain's intrinsic hash chain structure inherently tracks the provenance of its ordered blocks, the ordered transactions in each block, and the current state of the blockchain. In order to provide analysis of a specific historic transaction, all preceding transactions in all preceding blocks need to be run, which is time and resources intensive. Ruan et al. [8] proposed a system that collects the changed values of each blockchain address after each transaction execution and stores it in a Merkle tree. Contracts access this information using an efficient index, which can increase the contracts' applicability. In order to analyse contracts execution flow, additional data is needed. Devecsery et al. [10] discuss software systems that remember all operations, function calls, and states at any time and refer to them as *Eidetic software systems*. This historic
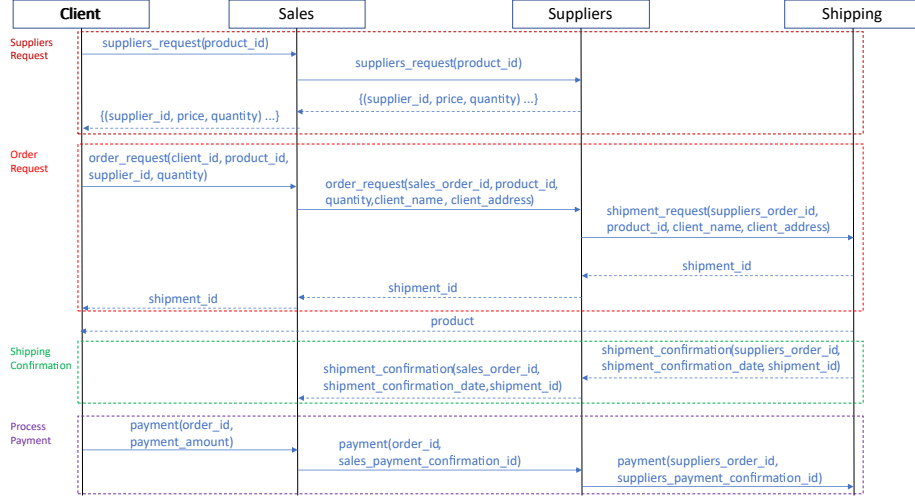
Fig. 1. Execution flow of a full order request session

information facilitates analysis of metadata, e.g., it enables queries that can answer questions about what states or calls affected other states or calls, and conversely what states or calls were affected by other states or calls. Inspired by this, we propose an *Eidetic blockchain system* that supports provenance for both data (blockchain states) and control flow (contracts). Our system captures the provenance of the contracts' execution flow (execution flow provenance), their parameters, and the relevant blockchain states before and after each contract call (Why-provenance). This provenance information can then be used to query the execution flow across time in different granularity levels, facilitate a more complete understanding of the contracts' execution environment, and ultimately assist in better security analysis of suspicious contracts behavior, as well as regulatory, quality, and maintenance management.

## II. SMART CONTRACT EXECUTION FLOW PROVENANCE

Existing approaches to capturing data provenance on blockchain [8] [18] gather information regarding the past history of blockchain states in a secure fashion. While, this is quite useful, for the purpose of finding defects in contracts code, and preventing potential blockchain abuse, it only provides information on what happened, rather than why or how it happened. Accurately tracking information flow in blockchain, along with its states and the contract execution flow that operates on the data, can help identify the root causes of anomalies, trace negative debugging effects, investigate malicious actors, and enable efficient audit tracking for regulation purposes. In the next sections we provide a motivating example followed by the description of our proposed system.

### A. Motivating Example

Alice, Bob, and Carol are friends and uPhone enthusiasts. Alice logs to an online retailer site to look for recent deals on the new phone model. She receives a list of suppliers that offer the phone alongside its price and quantity per supplier. Alice then picks the supplier that offers the best price and issues an order request. The Sales department receives the order request, and issues an order request of its own to

TABLE I
SUPPLIERS DATABASE INITIAL DATA

| supplier_id | product_id | product_quantity | product_price |
|---|---|---|---|
| 1 | 100 | 3 | 135 |
| 2 | 100 | 2 | 110 |

the Suppliers department. The Suppliers department issues a shipping request to the Shipping department for the product from the specific supplier to Alice. The Shipping department then issues a shipping identifier and returns it to the Suppliers department that returns it to Sales, that returns it to Alice. When the product is delivered, the Shipping department sends the shipping confirmation to Suppliers, that sends it to Sales. Sales charges Alice using her recorded payment information and issues a payment confirmation for Suppliers, which issues a payment confirmation for Shipping. The execution flow for Alice's example can be seen in Fig. 1. Alice then informs Bob about the supplier with the great deal she found. Bob logs to the site, issues a similar order request, and calls to update Carol, who logs to the site but can not find the deal. He calls customer service to inquire why.

### B. Entities, Contracts, and Data

Each department's operational logic can be encapsulated by a contract, i.e. Sales, Suppliers, and Shipping each has its own contract. Each department maintains its database state on the blockchain. Each department issues a contract call using their corresponding client, e.g., the Sales department uses a sales_client to call the Sales contract. The user that issues the order request uses its own client. A description of selected contract API functions that participate in the order_request scenario are provided in Table II. The initial Suppliers database can be seen in Table I. In our example we use two user clients to issue an order request: Alice with client_id 1 is the first to issue the order request for product_id 100 from supplier_id 2, which has only 2 units of the product. Bob with client_id 2 follows with an order request of the same product from the same supplier and orders an additional unit.

TABLE II
CONTRACTS DESCRIPTION

| Smart contract | API function | Description |
| --- | --- | --- |
| sales | order_request | Records the order details in the client_orders_DB and calls the order_request function in the Suppliers smart contract. The response should contain a shipment_id, which is updated in the order details. |
| suppliers | order_request | Records the order details in the sales_orders_DB and calls the shipping_request function in the Shipping smart contract. The response should contain a shipment_id, which is updated in the order details and is returned to the calling smart contract. |
| shipping | shipping_request | The client's physical address details are updated, a shipment_id is generated and is returned to the calling smart contract. |



Fig. 2. A session's contracts call graph

## III. THE PROPOSED SYSTEM

### A. Collecting provenance data

In our proposed system the contract execution infrastructure is modified to include a non-intrusive provenance collection mechanism that collects relevant contracts' execution information during runtime and is comprised of:

- Execution parameters, such as the caller (client or contract) name/address, the callee (contract) name/address, and the API function alongside its parameters to be invoked by the callee.
- Contracts call graph.
- Relevant states before and after the caller calls the callee.

When a caller calls the callee contract only the states that are read or written by the caller and callee are collected. This is done by modifying the blockchain context's get/set functions to record what addresses were read/written. The produced results provide a more concise view of the provenance and facilitate an easier analysis. The collected data are stored in a provenance database. A different process converts and exports the collected provenance data into a graph database such as neo4j [9] to enable visualizing and querying the provenance using the graph database query language.

Our system generates two types of provenance graphs. The first type is a contracts call graph, where new entity nodes are generated and connected per contract call to reflect the provenance of the call graph (Fig. 2). The second type is contracts parameters and state change graph, where for each contract call its parameters and the states before and after the call are captured and connected to the caller node (Fig. 3).

### B. Querying provenance data

A user client's session is comprised of multiple scenarios. Transactions that are related to a specific scenario share a unique scenario_id as well as a session_id that identifies the user client's session. The contracts call graph provides information on what contracts were called during each session and scenario and by whom (client/contract). Fig. 2 shows the contracts call graph of Alice's session. From this we can see whether all contracts were called in the correct sequence. Each node contains an entity id and name, the session id, the scenario id, and the contract API function that was invoked on the callee contract. To get more detailed information on each contract call, it is possible to further query the contracts' parameters and state change graph. Further, this graph can be used to answer queries regarding changes in different granularity levels, as the following examples explore:

*a) Querying changes across multiple sessions:* Returning to the motivating example, Alice and Carol both ordered the new uPhone from supplier id 2. Carol could not find the same supplier and called the customer service to inquire about it. A customer service representative queries the Suppliers database and sees that the supplier with id 2 has no more products. To reveal the history that led to this current state, the representative can issue a query that shows what order requests that involved supplier with id 2 were issued and their details. The result can be seen in table III. The table shows that customers with id 1 and 2 issued an order_request on the specific date and time for a product from supplier_id 2. The first line shows that after Alice issued the request, the product quantity changed into 1, and the second line shows that after Bob issued the request the product quantity changed into 0.

*b) Querying changes in a session across multiple scenarios:* Let's assume that Alice would like to know when her account was charged for the item. Table IV shows the query results for Alice's specific session in the Sales department's client_orders_DB. The first line shows that after the order_request the Sales department created a client order row with id 111. Since the shipping id is provided at the end of the order_request call it can be seen in the last column. The second line shows that the payment_confirmation_id was updated at the end of the process_payment call and its datetime.

9

| Contract call parameters | | | | API function parameters | | | | Suppliers db row after contract call | | |
|---|---|---|---|---|---|---|---|---|---|---|
| request datetime | API function | from | to | supplier id | client id | product id | product quantity | supplier id | product id | product quantity |
| 2020-01-10 2:02:34 | order_request | client | sales | 2 | 1 | 100 | 1 | 2 | 100 | 1 |
| 2020-01-10 3:16:15 | order_request | client | sales | 2 | 2 | 100 | 1 | 2 | 100 | 0 |

| Contract call parameters | | | | API function parameters | | | Client orders db row after contract call | | |
|---|---|---|---|---|---|---|---|---|---|
| request datetime | API function | from | to | sales order id | supplier id | payment confirmation id | sales order id | payment confirmation id | shipment id |
| 2020-01-10 2:02:34 | order_request | sales | suppliers | 111 | 2 | | 111 | | 555 |
| 2020-01-13 1:13:11 | process_payment | sales | suppliers | 111 | | 123 | 111 | 123 | 555 |



Fig. 3. Suppliers contract calls Shipping_request on Shipping contract

*c) Querying changes in a specific scenario and entity:*
Fig. 3 shows a drill down query result for the Suppliers node in the order_request scenario (marked by a red square as in Fig. 2), where the Suppliers department requests Shipping for the product ordered by Alice. The middle node represents the Suppliers department entity. The attached green node represents the parameters used in the contract call such as the contract and the API function names, and the request datetime; the attached pink node contains the function specific parameters such as Alice's name and address, and the suppliers order id; the blue node represents the relevant blockchain state before the Shipping contract call; and the yellow node the relevant blockchain state after the call. As can be seen, the top most nodes (marked by a black rectangle) are present only in the state after the call and represent the Suppliers' shipping order for Alice that was recorded by the Shipping contract.

## IV. CONCLUSION

The adoption of smart contract applications in managing valuable assets rapidly increases alongside the applications' size and complexity. As a result, potential for exploitation, the amount of coding errors and regulation requirements with hefty financial consequences increase as well. In order to facilitate root-cause analysis of anomalies and investigate defects or malicious activities in the smart contracts, the blockchain systems need to efficiently manage both data and execution flow. In this paper we propose a blockchain provenance collection and analysis system and explore how the system can non-intrusively capture relevant provenance information of the contracts' execution flow, their parameters, and the blockchain states before and after each contract call. We further explore how this information can be queried at different levels of granularity to facilitate better regulatory, quality, and maintenance management.

## REFERENCES

[1] Ethereum, https://ethereum.org/
[2] Hyperledger, https://www.hyperledger.org/
[3] Blockchain Applications For The Modern Nation, https://cryptodaily.co.uk/2018/03/blockchain-applications/
[4] Blockchain has grabbed the attention of investors, https://www.cnbc.com/2018/04/02/blockchain-has-grabbed-the-attention-of-investors.html
[5] Three Near-term Applications For Blockchain Technology, https://www.forbes.com/sites/forbesfinancecouncil/2018/03/28/three-near-term-applications-for-blockchain-technology/2/#6c9f49c6310d
[6] Glavic, B. and Alonso, G., 2009. Perm: Processing provenance and data on the same data model through query rewriting. ICDE. IEEE.
[7] Psallidas, F. and Wu, E., 2018. Smoke: Fine-grained lineage at interactive speed. VLDB.
[8] Ruan, P., Chen, G., Dinh, T.T.A., Lin, Q., Ooi, B.C. and Zhang, M., 2019. Fine-grained, secure and efficient data provenance on blockchain systems. VLDB.
[9] neo4j, https://neo4j.com/
[10] Devecsery, D., Chow, M., Dou, X., Flinn, J. and Chen, P.M., 2014. Eidetic systems. OSDI.
[11] Cheney, J., Chiticariu, L. and Tan, W.C., 2009. Provenance in databases: Why, how, and where. Foundations and Trends® in Databases.
[12] Miao, H., Chavan, A. and Deshpande, A., 2017, May. Provdb: Lifecycle management of collaborative analysis workflows. HILDA workshop. ACM.
[13] https://www.statista.com/statistics/1015362/worldwide-blockchain-technology-market-size/
[14] ETHNews: Hkg token has a bug and needs to be reissued (2017).
[15] Simonite, T.: $80 million hack shows the dangers of programmable money, June 2016. https://www.technologyreview.com
[16] Liang, X., Shetty, S., Tosh, D., Kamhoua, C., Kwiat, K. and Njilla, L., 2017, May. Provchain: A blockchain-based data provenance architecture in cloud environment with enhanced privacy and availability. CCGRID.
[17] Ramachandran, A. and Kantarcioglu, M., 2018, March. SmartProvenance: a distributed, blockchain based dataprovenance system. CODASPY. ACM.
[18] C. Xu, C. Zhang, and J. Xu. 2019. VChain: Enabling Verifiable Boolean Range Queries over Blockchain Databases. SIGMOD. ACM.