

End-to-End Formal Verification of Ethereum 2.0 Deposit Smart Contract (Solidity Reimplementation)

Daejun Park, Yi Zhang, and Grigore Rosu

Runtime Verification, Inc.

Abstract. We present our formal verification of the deposit smart contract, reimplemented in Solidity, whose correctness is critical for the security of Ethereum 2.0, a new Proof-of-Stake protocol for the Ethereum blockchain. The deposit contract employs an incremental Merkle tree algorithm whose correctness is highly nontrivial, and had not been proved before. We have formally verified the correctness of the compiled bytecode of the deposit contract to avoid the need to trust the underlying compiler. No critical issues were found in the Solidity implementation of the deposit contract.

1 Introduction

The deposit smart contract [4] is a gateway to join Ethereum 2.0 [5] that is a new sharded Proof-of-Stake (PoS) protocol which at its early stage, lives in parallel with the existing Proof-of-Work (PoW) chain, called Ethereum 1.x chain. Validators drive the entire PoS chain, called Beacon chain, of Ethereum 2.0. To be a validator, one needs to deposit a certain amount of Ether, as a “stake”, by sending a transaction (over the Ethereum 1.x network) to the deposit contract. The deposit contract records the history of deposits, and locks all the deposits in the Ethereum 1.x chain, which can be later claimed at the Beacon chain of Ethereum 2.0.¹ Note that the deposit contract is a one-way function; one can move her funds from Ethereum 1.x to Ethereum 2.0, but not vice versa.

The deposit contract, written in Solidity [8],² employs the Merkle tree [11] data structure to efficiently store the deposit history, where the tree is *dynamically* updated (i.e., leaf nodes are incrementally added in order from left to right) whenever a new deposit is received. The Merkle tree employed in this contract is very large: it has height 32, so it can store up to 2^{32} deposits. Since the size of the Merkle tree is huge, it is not practical to reconstruct the whole tree every time a new deposit is received.

To reduce both time and space complexity, thus saving the gas cost significantly, the contract implements an *incremental Merkle tree algorithm* [1]. The

¹ This deposit process will change at a later stage.

² Initially, the contract had been written in Vyper [9] but later it was reimplemented in Solidity due to the concerns [2] regarding the Vyper compiler. Indeed, we had also formally verified the previous Vyper implementation, which can be found at [17].

incremental algorithm enjoys $O(h)$ time and space complexity to reconstruct (more precisely, compute the root of) a Merkle tree of height h , while a naive algorithm would require $O(2^h)$ time or space complexity. The efficient incremental algorithm, however, leads to the deposit contract implementation being unintuitive, and makes it non-trivial to ensure its correctness. The correctness of the deposit contract, however, is critical for the security of Ethereum 2.0, since it is a gateway for becoming a validator. Considering the utmost importance of the deposit contract for the Ethereum blockchain, formal verification is demanded to ultimately guarantee its correctness.

In this paper, we present our formal verification of the deposit contract.³ The scope of verification is to ensure the correctness of the contract bytecode within a single transaction, without considering transaction-level or off-chain behaviors. We take the compiled bytecode as the verification target to avoid the need to trust the compiler.

We adopt a refinement-based verification approach. Specifically, our verification effort consists of the following two tasks:

- Verify that the incremental Merkle tree algorithm implemented in the deposit contract is *correct* w.r.t. the original full-construction algorithm.
- Verify that the compiled bytecode is *correctly generated* from the source code of the deposit contract.

Intuitively, the first task amounts to ensuring the correctness of the contract source code, while the second task amounts to ensuring the compiled bytecode being a sound refinement of the source code (i.e., translation validation of the compiler). This refinement-based approach allows us to avoid reasoning about the complex algorithmic details, especially specifying and verifying loop invariants, directly at the bytecode level. This separation of concerns helped us to save a significant amount of verification effort. See Section 2 for more details.

Verification Target. The specific target of our formal verification is the latest version (v0.12) of the deposit contract written in Solidity, provided that the contract is compiled by the Solidity compiler v0.6.8 with the optimization enabled (`-optimize-runs 5000000`).

2 Our Refinement-Based Verification Approach

We illustrate our refinement-based formal verification approach used in the deposit contract verification. We present our approach using the K framework and its verification infrastructure [18,20,13], but it can be applied to other program verification frameworks.

Let us consider a `sum` program that computes the summation from 1 to n :

```
int sum(int n) { int s = 0; int i = 1;
                while(i <= n) { s = s + i; i = i + 1; } return s; }
```

³ This was done as part of a contract funded by the Ethereum Foundation [6].

Given this program, we first manually write an abstract model of the program in the K framework [18]. Such a K model is essentially a state transition system of the program, and can be written as follows:

```
rule: sum(n) ⇒ loop(s: 0, i: 1, n: n)
rule: loop(s: s, i: i, n: n) ⇒ loop(s: s+i, i: i+1, n: n) when i ≤ n
rule: loop(s: s, i: i, n: n) ⇒ return(s) when i > n
```

These transition rules correspond to the initialization, the **while** loop, and the return statement, respectively. The indexed tuple $(s: s, i: i, n: n)$ represents the state of the program variables **s**, **i**, and **n**.⁴

Then, given the abstract model, we specify the functional correctness property in reachability logic [19], as follows:

```
claim: sum(n) ⇒ return( $\frac{n(n+1)}{2}$ ) when n > 0
```

This reachability claim says that $\text{sum}(n)$ will eventually return $\frac{n(n+1)}{2}$ in all possible execution paths, if n is positive. We verify this specification using the K reachability logic theorem prover [20], which requires us only to provide the following loop invariant:⁵

```
invariant: loop(s:  $\frac{i(i-1)}{2}$ , i: i, n: n) ⇒ return( $\frac{n(n+1)}{2}$ ) when 0 < i ≤ n + 1
```

Once we prove the desired property of the abstract model, we manually refine the model to a bytecode specification, by translating each transition rule of the abstract model into a reachability claim at the bytecode level, as follows:

```
claim: evm(pc: pcbegin, calldata: #bytes(32, n), stack: [], ...)
      ⇒ evm(pc: pcloophead, stack: [0, 1, n], ...)
claim: evm(pc: pcloophead, stack: [s, i, n], ...)
      ⇒ evm(pc: pcloophead, stack: [s+i, i+1, n], ...) when i ≤ n
claim: evm(pc: pcloophead, stack: [s, i, n], ...)
      ⇒ evm(pc: pcend, stack: [], output: #bytes(32, s), ...) when i > n
```

Here, the indexed tuple $\text{evm}(\text{pc}:_, \text{calldata}:_, \text{stack}:_, \text{output}:_)$ represents (part of) the Ethereum Virtual Machine (EVM) state, and $\#bytes(N, V)$ denotes a sequence of N bytes of the two's complement representation of V .

We verify this bytecode specification against the compiled bytecode using the same K reachability theorem prover [20,13]. Note that no loop invariant is needed in this bytecode verification, since each reachability claim involves only a bounded number of execution steps—specifically, the second claim involves only a single iteration of the loop.

Then, we manually prove the soundness of the refinement, which can be stated as follows: *for any EVM states σ_1 and σ_2 , if $\sigma_1 \Rightarrow \sigma_2$, then $\alpha(\sigma_1) \Rightarrow \alpha(\sigma_2)$* , where the abstraction function α is defined as follows:

⁴ Note that this abstract model can be also automatically derived by instantiating the language semantics with the particular program, if a formal semantics of the language is available (in the K framework).

⁵ The loop invariants in reachability logic mentioned here look different from those in Hoare logic. See the comparison between the two logic proof systems in [20, Section 4]. These loop invariants can be also seen as transition invariants [14].

2. The path can be computed by using only the left (i.e., node 3 and node 9) or right (i.e., node 14) sibling of each node in the path. All the other nodes are *not* needed for the path computation.
3. All the left siblings (i.e., node 3 and node 9) of the path are “finalized” in that they will never be updated in any subsequent execution of the algorithm. All the leaves that are a descendant of the finalized node are non-empty.
4. All the right siblings (i.e., node 14) are zero-hashes, that is, 0 for leaf nodes (at level 0), “hash(0,0)” for nodes at level 1, “hash(hash(0,0),hash(0,0))” for nodes at level 2, and so on. These zero-hashes are constant.

Now we describe the algorithm. To represent a Merkle tree of height h , the algorithm maintains only two arrays of length h , called **branch** and **zero_hashes** respectively, that store the left and right siblings of a path from a new leaf node to the root. When inserting a new data hash, the algorithm computes the path from the new leaf node to the root. Each node of the path can be computed in constant time, by retrieving only its left or right sibling from the **branch** or **zero_hashes** array. After the path computation, the **branch** array is updated to contain all the left siblings of a next new path that will be computed in the next run of the algorithm. Here the **branch** array update is done in constant time, since only a single element of the array needs to be updated, and the element has already been computed as part of the path computation.⁶ Note that the **zero_hashes** array is computed once at the very beginning when all the leaves are empty, and never be updated during the lifetime of the Merkle tree.

Complexity. Both the time and space complexity of the algorithm is linear in the tree height h . The space complexity is linear, because the size of the **branch** and **zero_hashes** arrays is h , and no other nodes are stored by the algorithm. The time complexity is also linear. For the path computation, the length of the path is h , and each node can be computed in constant time by using the two arrays. The **branch** array update can be also done in constant time as explained earlier.

Implementation and optimization. Figure 2 shows the pseudocode implementation of the incremental Merkle tree algorithm [1] that is employed in the deposit contract [4]. It consists of two main functions: **deposit** and **get_deposit_root**. The **deposit** function takes as input a new deposit hash, and inserts it into the Merkle tree. The **get_deposit_root** function computes and returns the root of the current partial Merkle tree whose leaves are filled with the deposit hashes received up to that point.

Specifically, the **deposit** function fills the first (leftmost) empty leaf node with a given deposit hash, and updates a single element of the **branch** array. The **get_deposit_root** function computes the tree root by traversing a path from the last (rightmost) non-empty leaf to the root.

As an optimization, the **deposit** function does not fully compute the path from the leaf to the root, but computes only a smaller partial path from the

⁶ See Appendix A for more details about updating the **branch** array.

```

1  # globals
2  zero_hashes: int[TREE_HEIGHT] = {0} # zero array
3  branch:      int[TREE_HEIGHT] = {0} # zero array
4  deposit_count: int = 0 # max: 2^TREE_HEIGHT - 1
5
6  fun constructor() -> unit:
7      i: int = 0
8      while i < TREE_HEIGHT - 1:
9          zero_hashes[i+1] = hash(zero_hashes[i], zero_hashes[i])
10         i += 1
11
12  fun deposit(value: int) -> unit:
13      assert deposit_count < 2^TREE_HEIGHT - 1
14      deposit_count += 1
15      size: int = deposit_count
16      i: int = 0
17      while i < TREE_HEIGHT:
18          if size % 2 == 1:
19              break
20          value = hash(branch[i], value)
21          size /= 2
22          i += 1
23      branch[i] = value
24
25  fun get_deposit_root() -> int:
26      root: int = 0
27      size: int = deposit_count
28      h: int = 0
29      while h < TREE_HEIGHT:
30          if size % 2 == 1: # size is odd
31              root = hash(branch[h], root)
32          else:             # size is even
33              root = hash(root, zero_hashes[h])
34          size /= 2
35          h += 1
36      return root

```

Fig. 2. Pseudocode implementation of the incremental Merkle tree algorithm employed in the deposit contract [4].

leaf to the node that is needed to update the **branch** array. Indeed, for all odd-numbered deposits (i.e., 1st deposit, 3rd deposit, \dots), such a partial path is empty, because the leaf node is the one needed for the **branch** array update. In that case, the **deposit** function returns immediately in constant time. For even-numbered deposits, the partial path is not empty but still much smaller than the full path in most cases. This optimization is useful when the tree root computation is not needed for every single partial Merkle tree. Indeed, in many cases, multiple deposit hashes are inserted at once, for which only the root of the last partial Merkle tree is needed.

Correctness. Consider a Merkle tree of height h employed in the deposit contract. Suppose that a sequence of **deposit** function calls are made, say **deposit**(v_1), **deposit**(v_2), \dots , and **deposit**(v_m), where $m < 2^h$. Then, the function call **get_deposit_root**() will return the root of the Merkle tree whose leaves are filled with the deposit data hashes v_1, v_2, \dots, v_m , respectively, in order from left to right, starting from the leftmost one.

Note that the correctness statement requires the condition $m < 2^h$, that is, the rightmost leaf must be kept empty, which means that the maximum number of deposits that can be stored in the tree using this incremental algorithm is $2^h - 1$ instead of 2^h .

The proof of the correctness is presented in Appendix A.

Remark. Since the **deposit** function reverts when **deposit_count** $\geq 2^{\text{TREE_HEIGHT}} - 1$, the loop in the **deposit** function cannot reach the last iteration, thus the loop bound (in line 17 of Figure 2) can be safely decreased to **TREE_HEIGHT** - 1.

4 Bytecode Verification of the Deposit Contract

Now we present the formal verification of the compiled bytecode of the deposit contract. The bytecode verification ensures that the compiled bytecode is a sound refinement of the source code. This rules out the need to trust the compiler.

As illustrated in Section 2, we first manually refined the abstract model (in which we proved the algorithm correctness) to the bytecode specification (Section 4.1). For the refinement, we consulted the ABI interface standard [3] (to identify, e.g., **calldata** and **output** in the illustrating example of Section 2), as well as the bytecode (to identify, e.g., the **pc** and **stack** information).⁷ Then, we used the KEVM verifier [13] to verify the compiled bytecode against the refined specification. We adopted the KEVM verifier to reason about all possible corner-case behaviors of the compiled bytecode, especially those introduced by certain unintuitive and questionable aspects of the underlying Ethereum Virtual Machine (EVM) [21]. This was possible because the KEVM verifier is derived

⁷ However, we want to note that the compiler can be augmented to extract such information, which can automate the refinement process to a certain extent. We leave that as future work.

from a complete formal semantics of the EVM, called KEVM [10]. Our formal specification and verification artifacts are publicly available at [16].

Let us elaborate on specific low-level behaviors verified against the bytecode. In addition to executing the incremental Merkle tree algorithm, most of the functions perform certain additional low-level tasks, and we verified that such tasks are correctly performed. Specifically, for example, given deposit data,⁸ the `deposit` function computes its 32-byte hash (called Merkleization) according to the SimpleSerialize (SSZ) specification [7]. The leaves of the Merkle tree store only the computed hashes instead of the original deposit data. The `deposit` function also emits a `DepositEvent` log that contains the original deposit data, where the log message needs to be encoded as a byte sequence following the contract event ABI specification [3]. Other low-level operations performed by those functions that we verified include: correct zero-padding for the 32-byte alignment, correct conversions from big-endian to little-endian, input bytes of the SHA2-256 hash function being correctly constructed, and return values being correctly serialized to byte sequences according to the ABI specification [3].

Our formal specification includes both positive and negative behaviors. The positive behaviors describe the desired behaviors of the contracts in a legitimate input state. The negative behaviors, on the other hand, describe how the contracts handle exceptional cases (e.g., when benign users feed invalid inputs by mistake, or malicious users feed crafted inputs to take advantage of the contracts). The negative behaviors are mostly related to security properties.

4.1 Summary of Bytecode Specification

We summarize the formal specification of the deposit contract bytecode that we verified. The full specification can be found at [15].

Constructor function `constructor` updates the storage as follows:

$$\text{zero_hashes}[i] \leftarrow ZH(i) \quad \text{for all } 1 \leq i < 32$$

where $ZH(i)$ is a 32-byte word that is recursively defined as follows:

$$\begin{aligned} ZH(i+1) &= \text{hash}(ZH(i) ++ ZH(i)) \quad \text{for } 0 \leq i < 31 \\ ZH(0) &= 0 \end{aligned}$$

where `hash` denotes the SHA2-256 hash function, and `++` denotes the byte concatenation.

Function `get_deposit_count` returns $LE_{64}(\text{deposit_count})$, where $LE_{64}(x)$ denotes the 64-bit little-endian representation of x (for $0 \leq x < 2^{64}$). That is, for a given $x = \sum_{0 \leq i < 8} (a_i \cdot 256^i)$, $LE_{64}(x) = \sum_{0 \leq i < 8} (a_{7-i} \cdot 256^i)$, where $0 \leq a_i < 256$. Note that $LE_{64}(\text{deposit_count})$ is always defined because of the contract invariant of `deposit_count` $< 2^{32}$. This function does not alter the storage state.

⁸ Each deposit data consists of the public key, the withdrawal credentials, the deposit amount, and the signature of the deposit owner.

Function *get_deposit_root* returns:

$$\text{hash}(RT(32) \mathbin{++} LE_{64}(\text{deposit_count}) \mathbin{++} 0_{[24]})$$

where $RT(32)$ is the Merkle tree root, recursively defined as follows:

$$RT(i+1) = \begin{cases} \text{hash}(\text{branch}[i] \mathbin{++} RT(i)), & \text{if } \lfloor \text{deposit_count}/2^i \rfloor \text{ is odd} \\ \text{hash}(RT(i) \mathbin{++} \text{zero_hashes}[i]), & \text{otherwise} \end{cases}$$

for $0 \leq i < 32$

$$RT(0) = 0$$

and $0_{[24]}$ denotes 24 zero-bytes. This function does not alter the storage state.

Function *deposit* updates the storage state as follows:

$$\begin{aligned} \text{deposit_count} &\leftarrow \text{old}(\text{deposit_count}) + 1 \\ \text{branch}[k] &\leftarrow ND(k) \end{aligned}$$

where $\text{old}(\text{deposit_count})$ denotes the value of `deposit_count` at the beginning of the function, k is the smallest integer less than 32 such that $\lfloor \frac{\text{old}(\text{deposit_count})+1}{2^k} \rfloor$ is odd,⁹ and $ND(K)$ is a 32-byte word that is recursively defined as follows:

$$ND(i+1) = \text{hash}(\text{branch}[i] \mathbin{++} ND(i)) \quad \text{for } 0 \leq i < 32$$

where $ND(0)$ denotes the deposit data root that is a Merkle proof of the deposit data that consists of the public key, the withdrawal credentials, the deposit amount, and the signature. The `deposit` function also emits a `DepositEvent` log that includes both the deposit data and the $\text{old}(\text{deposit_count})$ value. For the full details about the deposit data root computation and the `DepositEvent` log, refer to [15].

Negative behaviors. The contract reverts when either a call-value (i.e., `msg.value`) or a call-data (i.e., `msg.data`) is invalid. A call-value is invalid when it is non-zero but the called function is not payable (i.e., no `payable` annotation). A call-data is invalid when its size is less than 4 bytes, or its first four bytes do not match the signature of any public functions in the contract. Note that any extra contents in the call-data are silently ignored.¹⁰

The `deposit` function reverts if the tree is full, the deposit amount is less than the required minimum amount or not a multiple of Gwei (10^9 wei), or the call-data is not well-formed.

⁹ Note that such k always exists since we have $\text{old}(\text{deposit_count}) < 2^{32} - 1$ by the assertion at the beginning of the function.

¹⁰ We have not yet found an attack that can exploit this behavior.

5 Conclusion

We reported our end-to-end formal verification of the Ethereum 2.0 deposit contract. We adopted the refinement-based verification approach to ensure the end-to-end correctness of the contract while minimizing the verification effort. Specifically, we first proved that the incremental Merkle tree algorithm is correctly implemented in the contract, and then verified that the compiled bytecode is correctly generated from the source code. No critical bugs were found in the Solidity implementation of the deposit contract.¹¹ We conclude that the latest deposit contract (v0.12) will behave as expected—as specified in the formal specification [15], provided that the contract is compiled by the Solidity compiler v0.6.8 with the optimization enabled (`-optimize-runs 5000000`).

Trust base. The validity of the bytecode verification result assumes the correctness of the bytecode specification and the KEVM verifier. The algorithm correctness proof is partially mechanized—only the proof of major lemmas are mechanized in the K framework. The non-mechanized proofs are included in our trust base. The Solidity compiler is *not* in the trust base.

¹¹ Several issues found in the previous Vyper implementation have already been fixed in the Solidity reimplementation. See [17] for more details of the previous findings.

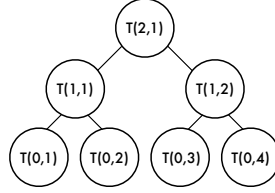
References

1. Buterin, V.: Progressive Merkle Tree. https://github.com/ethereum/research/blob/master/beacon_chain_impl/progressive_merkle_tree.py
2. ConsenSys Diligence: Vyper Security Review. <https://diligence.consys.net/audits/2019/10/vyper/>
3. Ethereum Foundation: Contract ABI Specification. <https://solidity.readthedocs.io/en/v0.6.1/abi-spec.html>
4. Ethereum Foundation: Ethereum 2.0 Deposit Contract. https://github.com/ethereum/eth2.0-specs/tree/master/deposit_contract
5. Ethereum Foundation: Ethereum 2.0 Specifications. <https://github.com/ethereum/eth2.0-specs>
6. Ethereum Foundation: Ethereum Foundation Spring 2019 Update. <https://blog.ethereum.org/2019/05/21/ethereum-foundation-spring-2019-update/>
7. Ethereum Foundation: SimpleSerialize (SSZ). <https://github.com/ethereum/eth2.0-specs/tree/master/ssz>
8. Ethereum Foundation: Solidity. <https://solidity.readthedocs.io>
9. Ethereum Foundation: Vyper. <https://vyper.readthedocs.io>
10. Hildenbrandt, E., Saxena, M., Zhu, X., Rodrigues, N., Daian, P., Guth, D., Moore, B., Zhang, Y., Park, D., Ștefănescu, A., Roșu, G.: Kevm: A complete semantics of the ethereum virtual machine. In: Proceedings of the 31st IEEE Computer Security Foundations Symposium. CSF 2018 (2018)
11. Merkle, R.C.: A digital signature based on a conventional encryption function. In: A Conference on the Theory and Applications of Cryptographic Techniques on Advances in Cryptology. CRYPTO '87 (1988)
12. NIST: Perfect Binary Tree. <https://xlinux.nist.gov/dads/HTML/perfectBinaryTree.html>
13. Park, D., Zhang, Y., Saxena, M., Daian, P., Roșu, G.: A Formal Verification Tool for Ethereum VM Bytecode. In: Proceedings of the 26th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering. ESEC/FSE 2018 (2018)
14. Podelski, A., Rybalchenko, A.: Transition invariants. In: Proceedings of the 19th Annual IEEE Symposium on Logic in Computer Science. LICS 2004 (2004)
15. Runtime Verification, Inc.: Formal Verification of Deposit Contract Bytecode. <https://github.com/runtimeverification/deposit-contract-verification/tree/master/bytecode-verification>
16. Runtime Verification, Inc.: Formal Verification of Ethereum 2.0 Deposit Contract. <https://github.com/runtimeverification/deposit-contract-verification>
17. Runtime Verification, Inc.: Formal Verification of Ethereum 2.0 Deposit Contract (Vyper Implementation). <https://github.com/runtimeverification/verified-smart-contracts/tree/master/deposit>
18. Serbanuta, T., Arusoiaie, A., Lazar, D., Ellison, C., Lucanu, D., Rosu, G.: The K primer (version 3.3). Electr. Notes Theor. Comput. Sci. **304**, 57–80 (2014)
19. Ștefănescu, A., Ciobaca, S., Mereuta, R., Moore, B.M., Serbanuta, T., Rosu, G.: All-Path Reachability Logic. Logical Methods in Computer Science **15**(2) (2019)
20. Ștefănescu, A., Park, D., Yuwen, S., Li, Y., Rosu, G.: Semantics-Based Program Verifiers for All Languages. In: Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications. OOPSLA 2016 (2016)
21. Wood, G.: Ethereum: A Secure Decentralised Generalised Transaction Ledger. <https://ethereum.github.io/yellowpaper/paper.pdf>

A Formalization and Correctness Proof of the Incremental Merkle Tree Algorithm

We formalize the incremental Merkle tree algorithm [1], especially the one employed in the deposit contract [4], and prove its correctness w.r.t. the original full-construction Merkle tree algorithm [11].

Notations. Let T be a perfect binary tree [12] (i.e., every node has exactly two child nodes) of height h , and $T(l, i)$ denote its node at level l and index i , where the level of leaves is 0, and the index of the left-most node is 1. For example, if $h = 2$, then $T(2, 1)$ denotes the root whose children are $T(1, 1)$ and $T(1, 2)$, and the leaves are denoted by $T(0, 1)$, $T(0, 2)$, $T(0, 3)$, and $T(0, 4)$, as follows:



We write $\llbracket T(l, i) \rrbracket$ to denote the value of the node $T(l, i)$, but we omit $\llbracket \cdot \rrbracket$ when the meaning is clear in the context.

Let us define two functions, \uparrow and \uparrow^k , as follows:

$$\uparrow x = \lceil x/2 \rceil \quad (1)$$

$$\uparrow^k x = \lfloor x/2 \rfloor \quad (2)$$

Moreover, let us define $\uparrow^k x = \uparrow(\uparrow^{k-1} x)$ for $k \geq 2$, $\uparrow^1 x = \uparrow x$, and $\uparrow^0 x = x$. Let $\{T(k, \uparrow^k x)\}_{k=0}^h$ be a path $\{T(0, \uparrow^0 x), T(1, \uparrow^1 x), T(2, \uparrow^2 x), \dots, T(h, \uparrow^h x)\}$. We write $\{T(k, \uparrow^k x)\}_k$ if h is clear in the context. Let us define \uparrow^k and $\{T(k, \uparrow^k x)\}_k$ similarly. For the presentation purpose, let $T(l, 0)$ denote a dummy node which has the parent $T(l+1, 0)$ and the children $T(l-1, 0)$ and $T(l-1, 1)$. Note that, however, these dummy nodes are only conceptual, allowing the aforementioned paths to be well-defined, but *not* part of the tree at all.

In this notation, for a non-leaf, non-root node of index i , its left child index is $2i-1$, its right child index is $2i$, and its parent index is $\uparrow i$. Also, note that $\{T(k, \uparrow^k m)\}_k$ is the path starting from the m -th leaf going all the way up to the root.

First, we show that two paths $\{T(k, \uparrow^k x)\}_k$ and $\{T(k, \uparrow^k (x-1))\}_k$ are parallel with a “distance” of 1.

Lemma 1. *For all $x \geq 1$, and $k \geq 0$, we have:*

$$(\uparrow^k x) - 1 = \uparrow^k (x - 1) \quad (3)$$

Proof. Let us prove by induction on k . When $k = 0$, we have $(\uparrow^0 x) - 1 = x - 1 = \uparrow^0 (x - 1)$. When $k = 1$, we have two cases:

– When x is odd, that is, $x = 2y + 1$ for some $y \geq 0$:

$$(\uparrow x) - 1 = (\uparrow (2y + 1)) - 1 = \left\lceil \frac{2y + 1}{2} \right\rceil - 1 = y = \left\lfloor \frac{2y}{2} \right\rfloor = \uparrow 2y = \uparrow (x - 1)$$

– When x is even, that is, $x = 2y$ for some $y \geq 1$:

$$(\uparrow x) - 1 = (\uparrow 2y) - 1 = \left\lceil \frac{2y}{2} \right\rceil - 1 = y - 1 = \left\lfloor \frac{2y - 1}{2} \right\rfloor = \uparrow (2y - 1) = \uparrow (x - 1)$$

Thus, we have:

$$(\uparrow x) - 1 = \uparrow (x - 1) \quad (4)$$

Now, assume that (3) holds for some $k = l \geq 1$. Then,

$$\begin{aligned} \uparrow^{l+1} x &= \uparrow (\uparrow^l x) && \text{(By the definition of } \uparrow^k) \\ &= \uparrow ((\uparrow^l (x - 1)) + 1) && \text{(By the assumption)} \\ &= (\uparrow (\uparrow^l (x - 1))) + 1 && \text{(By Equation 4)} \\ &= \uparrow^{l+1} (x - 1) + 1 && \text{(By the definition of } \uparrow^k) \end{aligned}$$

which concludes.

Now let us define the Merkle tree.

Definition 1. A perfect binary tree T of height h is a Merkle tree [11], if the leaf node contains data, and the non-leaf node's value is the hash of its children's, i.e.,

$$\forall 0 < l \leq h. \forall 0 < i \leq 2^{h-l}. T(l, i) = \text{hash}(T(l - 1, 2i - 1), T(l - 1, 2i)) \quad (5)$$

Let T_m be a partial Merkle tree up-to m whose first m leaves contain data and the other leaves are zero, i.e.,

$$T_m(0, i) = 0 \quad \text{for all } m < i \leq 2^h \quad (6)$$

Let Z be the zero Merkle tree whose leaves are all zero, i.e., $Z(0, i) = 0$ for all $0 < i \leq 2^h$. That is, $Z = T_0$. Since all nodes at the same level have the same value in Z , we write $Z(l)$ to denote the value at the level l , i.e., $Z(l) = Z(l, i)$ for any $0 < i \leq 2^{h-l}$.

Now we formulate the relationship between the partial Merkle trees. Given two partial Merkle trees T_{m-1} and T_m , if their leaves agree up-to $m - 1$, then they only differ on the path $\{T_m(k, \uparrow^k m)\}_k$. This is formalized in Lemma 2.

Lemma 2. Let T_m be a partial Merkle tree up-to $m > 0$ of height h , and let T_{m-1} be another partial Merkle tree up-to $m - 1$ of the same height. Suppose their leaves agree up to $m - 1$, that is, $T_{m-1}(0, i) = T_m(0, i)$ for all $1 \leq i \leq m - 1$. Then, for all $0 \leq l \leq h$, and $1 \leq i \leq 2^{h-l}$,

$$T_{m-1}(l, i) = T_m(l, i) \quad \text{when } i \neq \uparrow^l m \quad (7)$$

Proof. Let us prove by induction on l . When $l = 0$, we immediately have $T_{m-1}(0, i) = T_m(0, i)$ for any $i \neq m$ by the premise and Equation 6. Now, assume that (7) holds for some $l = k$. Then by Equation 5, we have $T_{m-1}(k+1, i) = T_m(k+1, i)$ for any $i \neq \uparrow^k m = \uparrow^{k+1} m$, which concludes.

Corollary 1 induces a *linear-time* incremental Merkle tree insertion algorithm [1].

Corollary 1. T_m can be constructed from T_{m-1} by computing only $\{T_m(k, \uparrow^k m)\}_k$, the path from the new leaf, $T_m(0, m)$, to the root.

Proof. By Lemma 2.

Let us formulate more properties of partial Merkle trees.

Lemma 3. Let T_m be a partial Merkle tree up-to m of height h , and Z be the zero Merkle tree of the same height. Then, for all $0 \leq l \leq h$, and $1 \leq i \leq 2^{h-l}$,

$$T_m(l, i) = Z(l) \quad \text{when } i > \uparrow^l m \quad (8)$$

Proof. Let us prove by induction on l . When $l = 0$, we immediately have $T_m(0, i) = Z(0) = 0$ for any $m < i \leq 2^h$ by Equation 6. Now, assume that (8) holds for some $0 \leq l = k < h$. First, for any $i \geq (\uparrow^{k+1} m) + 1$, we have:

$$2i - 1 \geq (2 \uparrow^{k+1} m) + 1 = 2 \left\lceil \frac{\uparrow^k m}{2} \right\rceil + 1 \geq 2 \frac{\uparrow^k m}{2} + 1 = (\uparrow^k m) + 1 \quad (9)$$

Then, for any $\uparrow^{k+1} m < i \leq 2^{h-(k+1)}$, we have:

$$\begin{aligned} T_m(k+1, i) &= \text{hash}(T_m(k, 2i-1), T_m(k, 2i)) && \text{(By Equation 5)} \\ &= \text{hash}(Z(k), Z(k)) && \text{(By Equations 8 and 9)} \\ &= Z(k+1) && \text{(By the definition of } Z) \end{aligned}$$

which concludes.

Lemma 4 induces a *linear-space* incremental Merkle tree insertion algorithm.

Lemma 4. A path $\{T_m(k, \uparrow^k m)\}_k$ can be computed by using only two other paths, $\{T_{m-1}(k, \uparrow^k (m-1))\}_k$ and $\{Z(k)\}_k$.

Proof. We will construct the path from the leaf, $T_m(0, m)$, which is given. Suppose we have constructed the path up to $T_m(q, \uparrow^q m)$ for some $q > 0$ by using only two other sub-paths, $\{T_{m-1}(k, \uparrow^k (m-1))\}_{k=0}^{q-1}$ and $\{Z(k)\}_{k=0}^{q-1}$. Then, to construct $T_m(q+1, \uparrow^{q+1} m)$, we need the sibling of $T_m(q, \uparrow^q m)$, where we have two cases:

- Case $(\uparrow^q m)$ is odd. Then, we need the right-sibling $T_m(q, (\uparrow^q m) + 1)$, which is $Z(q)$ by Lemma 3.

- Case $(\uparrow^q m)$ is even. Then, we need the left-sibling $T_m(q, (\uparrow^q m) - 1)$, which is $T_m(q, \uparrow^q(m - 1))$ by Lemma 1, which is in turn $T_{m-1}(q, \uparrow^q(m - 1))$ by Lemma 2.

By the mathematical induction on k , we conclude.

Lemma 5. *Let $h = \text{TREE_HEIGHT}$. For any integer $0 \leq m < 2^h$, the two paths $\{T_m(k, \uparrow^k m)\}_k$ and $\{T_{m+1}(k, \uparrow^k(m+1))\}_k$ always converge, that is, there exists unique $0 \leq l \leq h$ such that:*

$$(\uparrow^k m) + 1 = \uparrow^k(m + 1) \text{ is even for all } 0 \leq k < l \quad (10)$$

$$(\uparrow^k m) + 1 = \uparrow^k(m + 1) \text{ is odd for } k = l \quad (11)$$

$$\uparrow^k m = \uparrow^k(m + 1) \text{ for all } l < k \leq h \quad (12)$$

$$T_m(k, \uparrow^k m) = T_{m+1}(k, \uparrow^k(m + 1)) \text{ for all } l < k \leq h \quad (13)$$

Proof. Equation 12 follows from Equation 11, since for an odd integer x , $\uparrow(x - 1) = \uparrow x$. Also, Equation 13 follows from Lemma 2, since $\uparrow^k(m + 1) = (\uparrow^k m) + 1 \neq \uparrow^k m = \uparrow^k(m + 1)$ by Lemma 1 and Equation 12. Thus, we only need to prove the unique existence of l satisfying (10) and (11). The existence of l is obvious since $1 \leq m + 1 \leq 2^h$, and one can find the smallest l satisfying (10) and (11). Now, suppose there exist two different $l_1 < l_2$ satisfying (10) and (11). Then, $\uparrow^{l_1}(m + 1)$ is odd since l_1 satisfies (11), while $\uparrow^{l_1}(m + 1)$ is even since l_2 satisfies (10), which is a contradiction, thus l is unique, and we conclude.

A.1 Pseudocode

Figure 2 shows the pseudocode of the incremental Merkle tree algorithm [1] that is employed in the deposit contract [4]. It maintains a global counter `deposit_count` to keep track of the number of deposits made, and two global arrays `zero_hashes` and `branch`, which corresponds to Z (Definition 1) and a certain part of $\{T_m(k, \uparrow^k m)\}_k$, where m denotes the value of `deposit_count`. The `constructor` function is called once at the beginning to initialize `zero_hashes` which is never updated later. The `deposit` function inserts a given new leaf value in the tree by incrementing `deposit_count` and updating only a single element of `branch`. The `get_deposit_root` function computes the root of the current partial Merkle tree T_m .

Since the loops are bounded to the tree height and the size of global arrays is equal to the tree height, it is clear that both time and space complexities of the algorithm are linear.

A.2 Correctness Proof

Now we prove the correctness of the incremental Merkle tree algorithm shown in Figure 2.

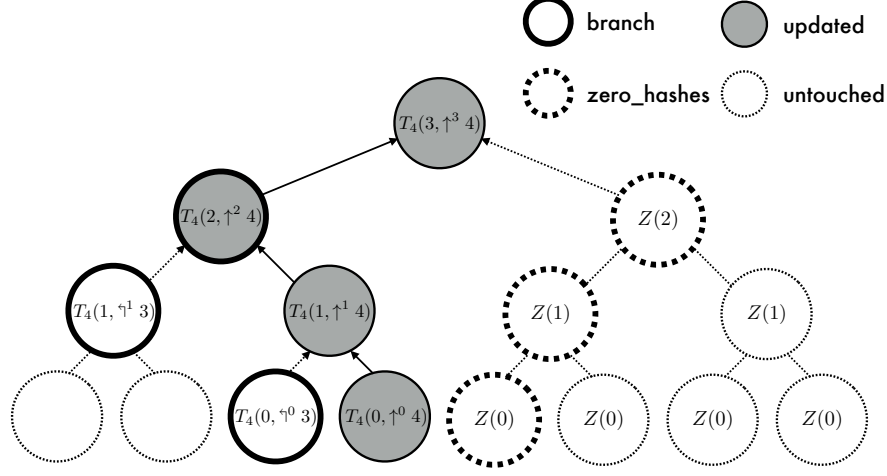


Fig. 3. A partial Merkle tree T_4 of height 3, illustrating the incremental Merkle tree algorithm shown in Figure 2, where $\text{TREE_HEIGHT} = 3$. The bold-lined nodes correspond to the **branch** array. The bold-dotted-lined nodes correspond to the **zero_hashes** array. The **get_deposit_root** function computes the gray nodes by using only the bold-lined nodes (i.e., **branch**) and the bold-dotted-lined nodes (i.e., **zero_hashes**), where $\text{deposit_count} = 4$.

Theorem 1 (Correctness of Incremental Merkle Tree Algorithm). *Suppose that the **constructor** function is executed at the beginning, followed by a sequence of **deposit** function calls, say $\text{deposit}(v_1)$, $\text{deposit}(v_2)$, \dots , and $\text{deposit}(v_m)$, where $m < 2^{\text{TREE_HEIGHT}}$. Then, the function call **get_deposit_root**() will return the root of the partial Merkle tree T_m such that $T_m(0, i) = v_i$ for all $1 \leq i \leq m$.*

Proof. By Lemmas 6, 7, 8, and 9.

Note that the correctness theorem requires the condition $m < 2^h$, where h is the tree height, that is, the rightmost leaf must be kept empty, which means that the maximum number of deposits that can be stored in the tree using this incremental algorithm is $2^h - 1$ instead of 2^h .

Lemma 6 (init). *Once **init** is executed, **zero_hashes** denotes Z , that is,*

$$\text{zero_hashes}[k] = Z(k) \quad (14)$$

for $0 \leq k < \text{TREE_HEIGHT}$.

Proof. By the implementation of **init** and the definition of Z in Definition 1.

Lemma 7 (deposit). *Suppose that, before executing **deposit**, we have:*

$$\text{deposit_count} = m < 2^{\text{TREE_HEIGHT}} - 1 \quad (15)$$

$$\text{branch}[k] = T_m(k, \uparrow^k m) \quad \text{if } \uparrow^k m \text{ is odd} \quad (16)$$

Then, after executing **deposit**(v), we have:

$$\mathbf{deposit_count}' = m + 1 \leq 2^{\mathbf{TREE_HEIGHT}} - 1 \quad (17)$$

$$\mathbf{branch}'[k] = T_{m+1}(k, \wr^k(m+1)) \quad \text{if } \wr^k(m+1) \text{ is odd} \quad (18)$$

for any $0 \leq k < \mathbf{TREE_HEIGHT}$, where:

$$T_{m+1}(0, m+1) = v \quad (19)$$

Proof. Let $h = \mathbf{TREE_HEIGHT}$. Equation 17 is obvious by the implementation of **deposit**. Let us prove Equation 18. Let l be the unique integer described in Lemma 5. We claim that **deposit** updates only **branch**[l] to be $T_{m+1}(l, \wr^l(m+1))$. Then, for all $0 \leq k < l$, $\wr^k(m+1)$ is not odd. For $k = l$, we conclude by the aforementioned claim. For $l < k \leq h$, we conclude by Equation 13 and the fact that **branch**[k] is not modified (by the aforementioned claim).

Now, let us prove the aforementioned claim. Since **branch** is updated only at line 23, we only need to prove $i = l$ and **value** = $T_{m+1}(l, \wr^l(m+1))$ at that point. We claim the following loop invariant at line 17:

$$i = i < \mathbf{TREE_HEIGHT} \quad (20)$$

$$\mathbf{value} = T_{m+1}(i, \wr^i(m+1)) \quad (21)$$

$$\mathbf{size} = \wr^i(m+1) \quad (22)$$

$$\wr^k(m+1) \text{ is even for any } 0 \leq k < i \quad (23)$$

Note that i cannot reach $\mathbf{TREE_HEIGHT}$, since $(m+1) < 2^{\mathbf{TREE_HEIGHT}}$. Thus, by the loop invariant, we have the following after the loop at line 23:

$$i = i < \mathbf{TREE_HEIGHT} \quad (24)$$

$$\mathbf{value} = T_{m+1}(i, \wr^i(m+1)) \quad (25)$$

$$\mathbf{size} = \wr^i(m+1) \text{ is odd} \quad (26)$$

$$\wr^k(m+1) \text{ is even for any } 0 \leq k < i \quad (27)$$

Moreover, by Lemma 5, we have $i = l$, which suffices to conclude the aforementioned claim.

Now we only need to prove the loop invariant. First, at the beginning of the first iteration, we have $i = 0$, **value** = $v = T_{m+1}(0, m+1)$ by (19), and **size** = $(m+1)$, which satisfies the loop invariant. Now, assume that the invariant holds at the beginning of the i^{th} iteration that does not reach the **break** statement at line 19 (i.e., **size** = $\wr^i(m+1)$ is even). Then, $i' = i + 1$, **size**' = $\wr^{i+1}(m+1)$, and:

$$\begin{aligned} T_{m+1}(i+1, \wr^{i+1}(m+1)) &= \text{hash}(T_{m+1}(i, \wr^i m), T_{m+1}(i, \wr^i(m+1))) \\ &\quad \text{(by Equation 10)} \\ &= \text{hash}(T_m(i, \wr^i m), \mathbf{value}) \\ &\quad \text{(by Lemmas 1 \& 2 and Equation 21)} \\ &= \text{hash}(\mathbf{branch}[i], \mathbf{value}) \quad \text{(by Equations 16 \& 10)} \\ &= \mathbf{value}' \end{aligned}$$

Thus, the loop invariant holds at the beginning of the $(i + 1)^{\text{th}}$ iteration as well, and we conclude.

Lemma 8 (Contract Invariant). *Let $m = \text{deposit_count}$. Then, once `init` is executed, the following contract invariant holds. For all $0 \leq k < \text{TREE_HEIGHT}$,*

1. $\text{zero_hashes}[k] = Z(k)$
2. $\text{branch}[k] = T_m(k, \uparrow^k m)$ if $\uparrow^k m$ is odd
3. $\text{deposit_count} \leq 2^{\text{TREE_HEIGHT}} - 1$

Proof. Let us prove each invariant item.

1. By Lemma 6, and the fact that `zero_hashes` is updated by only `init`.
2. By Lemma 7, and the fact that `branch` is updated by only `deposit`.
3. By the assertion of `deposit` (at line 13 of Figure 2), and the fact that `deposit_count` is updated by only `deposit`.

Lemma 9 (`get_deposit_root`). *The `get_deposit_root` function computes the path $\{T_m(k, \uparrow^k(m+1))\}_k$ and returns the root $T_m(h, 1)$, given a Merkle tree T_m of height h , that is, $\text{deposit_count} = m < 2^h$ and $\text{TREE_HEIGHT} = h$ when `get_deposit_root` is invoked.*

Proof. We claim the following loop invariant at line 29, which suffices to conclude the main claim.

$$\begin{aligned} \mathbf{h} &= k \quad \text{where } 0 \leq k \leq h \\ \mathbf{size} &= \uparrow^k m \\ \mathbf{root} &= T_m(k, \uparrow^k(m+1)) \end{aligned}$$

Now let us prove the above loop invariant claim by the mathematical induction on k . The base case ($k = 0$) is trivial, since $\uparrow^0 m = m$, $\uparrow^0(m+1) = m+1$, and $T_m(0, m+1) = 0$ by Definition 1. Assume that the loop invariant holds for some $k = l$. Let \mathbf{h}' , \mathbf{size}' , and \mathbf{root}' denote the values at the next iteration $k = l + 1$. Obviously, we have $\mathbf{h}' = l + 1$ and $\mathbf{size}' = \uparrow^{l+1} m$. Also, we have $(\uparrow^l m) + 1 = \uparrow^l(m+1)$ by Lemma 1. Now, we have two cases:

- Case $\mathbf{size} = \uparrow^l m$ is odd. Then, $\uparrow^l(m+1)$ is even. Thus,

$$\begin{aligned} T_m(l+1, \uparrow^{l+1}(m+1)) &= \text{hash}(T_m(l, \uparrow^l m), T_m(l, \uparrow^l(m+1))) \\ &= \text{hash}(\text{branch}[l], \mathbf{root}) \quad (\text{by Lemma 8}) \\ &= \mathbf{root}' \end{aligned}$$

- Case $\mathbf{size} = \uparrow^l m$ is even. Then, $\uparrow^l(m+1)$ is odd. Thus,

$$\begin{aligned} T_m(l+1, \uparrow^{l+1}(m+1)) &= \text{hash}(T_m(l, \uparrow^l(m+1)), T_m(l, (\uparrow^l(m+1)) + 1)) \\ &= \text{hash}(\mathbf{root}, Z(l)) \quad (\text{by Lemma 3}) \\ &= \text{hash}(\mathbf{root}, \text{zero_hashes}[l]) \quad (\text{by Lemma 8}) \\ &= \mathbf{root}' \end{aligned}$$

Thus, we have $\mathbf{root}' = T_m(l+1, \uparrow^{l+1}(m+1))$, which concludes.

Mechanized Proofs. The loop invariant proofs of Lemma 7 and Lemma 9 are mechanized in the K framework, which can be found at [16].