

Formal Verification of Incremental Merkle Tree Algorithm

Daejun Park and Yi Zhang

Runtime Verification, Inc.

Abstract. We formalize the incremental Merkle tree algorithm [1,5], and prove its correctness w.r.t. the original one [3].

1 Introduction

The deposit contract of Eth 2.0 [1], written in Vyper, is a smart contract that records a set of deposits for the Beacon chain [2]. The essence of the contract is the implementation of a Merkle tree that stores the set of deposits, where the tree is updated whenever a new deposit is received. The Merkle tree employed in the contract is expected to be very large. Indeed, a Merkle tree of height 32, which can store up to 2^{32} deposits, is implemented in the current version¹ of the contract. 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 requirement (thus saving the gas cost), the contract implements the incremental Merkle tree algorithm [5]. The 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. Specifically, the algorithm maintains two arrays of length h , and each construction of a new tree requires to compute only a chain starting from the new leaf (i.e., the new deposit) to the root, where the computation of the chain requires only the two arrays, thus achieving the linear time and space complexity in the height of a tree. The efficient incremental algorithm, however, leads to the deposit contract implementation being unintuitive, and makes it non-trivial to ensure its correctness. Considering the utmost importance of the deposit contract, formal verification is demanded (indeed, the only known way) to ultimately guarantee the correctness of the contract.

2 Formalization of Incremental Merkle Tree Algorithm

Notations Let T be a perfect binary tree [4] (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 leafs is 0, and the index of the left-most node is 1. For example,

¹ https://github.com/ethereum/deposit_contract/blob/ed9c81dd6788142d22106df93d5654578063eb32/deposit_contract/contracts/validator_registration.v.py

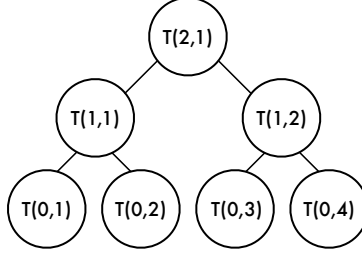


Fig. 1. A Merkle tree of height 2. We write $T(l, i)$ to denote the node of a tree T at the level l and the index i , where the level of leafs is 0, and the index of the left-most node is 1.

if $h = 2$, then $T(2, 1)$ denotes the root whose children are $T(1, 1)$ and $T(1, 2)$, and the leafs are denoted by $T(0, 1)$, $T(0, 2)$, $T(0, 3)$, and $T(0, 4)$, as shown in Figure 1. 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 \downarrow , as follows:

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

$$\downarrow x = \lfloor x/2 \rfloor \quad (2)$$

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 chain $\{T(0, \uparrow^0 x), T(1, \uparrow^1 x), T(2, \uparrow^2 x), \dots, T(h, \uparrow^h x)\}$. We write $\{T(k, \uparrow^k x)\}_k$ if h is clear in the context. Let us define \downarrow^k and $\{T(k, \downarrow^k x)\}_k$ 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 chains 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 chain starting from the m -th leaf going all the way up to the root.

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

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

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

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

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

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

– 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,

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

which concludes.

Now let us define the Merkle tree.

Definition 1. A perfect binary tree T of height h is a Merkle tree [3], 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 leafs contain data and the other leafs 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 leafs 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.

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 leafs 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 the equation (6). Now, assume that (7) holds for some $l = k$. Then by the equation 5, we have $T_{m-1}(k + 1, i) = T_m(k + 1, i)$ for any $i \neq \uparrow (\uparrow^k m) = \uparrow^{k+1} m$, which concludes.

Lemma 2 induces an incremental Merkle tree insertion algorithm [5].

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

Proof. By Lemma 2.

Let us formulate more properties of a partial Merkle tree.

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 the 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:

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

which concludes.

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

Proof. We will construct the chain from the leaf, $T_m(0, m)$, which is given. Suppose we have constructed the chain up to $T_m(q, \uparrow^q m)$ for some $q > 0$ by using only two other sub-chains, $\{T_{m-1}(k, \uparrow^k (m-1))\}_{k=0}^{q-1}$ 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 chains $\{T_m(k, \uparrow^k m)\}_k$ 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. The equation 12 follows from the equation 11, since for an odd integer x , $\uparrow(x-1) = \uparrow x$. Also, the 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 the 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 contradiction, thus l is unique, and we conclude.

3 Verification of Incremental Merkle Tree Algorithm

Now we verify the correctness of the incremental Merkle tree algorithm given in Figure 2.

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

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

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:

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

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

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

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

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

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

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

```

1  # globals
2  zerohashes: 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 init() -> unit:
7      i: int = 0
8      while i < TREE_HEIGHT - 1:
9          zerohashes[i+1] = hash(zerohashes[i], zerohashes[i])
10         i += 1
11
12  fun deposit(value: int) -> unit:
13      assert deposit_count < 2^TREE_HEIGHT - 1
14      deposit_count += 1
15
16      size: int = deposit_count
17      i: int = 0
18      while i < TREE_HEIGHT:
19          if size % 2 == 1:
20              break
21          size /= 2
22          i += 1
23
24      j: int = 0
25      while j < i:
26          value = hash(branch[j], value)
27          j += 1
28      branch[i] = value
29
30  fun get_deposit_root() -> int:
31      root: int = 0
32      size: int = deposit_count
33      h: int = 0
34      while h < TREE_HEIGHT:
35          if size % 2 == 1: # size is odd
36              root = hash(branch[h], root)
37          else:             # size is even
38              root = hash(root, zerohashes[h])
39          size /= 2
40          h += 1
41      return root

```

Fig. 2. Pseudocode implementation of the incremental Merkle tree algorithm employed in the deposit contract [1].

Proof. Let $h = \text{TREE_HEIGHT}$. The equation 17 is obvious by the implementation of **deposit**. Let us prove the equation 18. Let l be the unique integer described in Lemma 5. We claim that **deposit** updates only **branch**[l] to be $T_{m+1}(l, \gamma^l(m+1))$. Then, for all $0 \leq k < l$, $\gamma^k(m+1)$ is not odd. For $k = l$, we conclude by the aforementioned claim. For $l < k \leq h$, we conclude by the 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 28, we only need to prove $i = l$ and **value** = $T_{m+1}(l, \gamma^l(m+1))$ at that point. We claim the following loop invariants of the two inner loops in **deposit**.

– First loop invariant at line 18.

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

$$\text{size} = \gamma^i(m+1) \quad (21)$$

$$\gamma^k(m+1) \text{ is even for any } 0 \leq k < i \quad (22)$$

Note that i cannot reach **TREE_HEIGHT**, since $(m+1) < 2^{\text{TREE_HEIGHT}}$. It is trivial to show the above loop invariant. By the loop invariant, we have the following after the loop at line 23:

$$i = i < \text{TREE_HEIGHT} \quad (23)$$

$$\text{size} = \gamma^i(m+1) \text{ is odd} \quad (24)$$

$$\gamma^k(m+1) \text{ is even for any } 0 \leq k < i \quad (25)$$

Moreover, by Lemma 5, we have $i = l$.

– Second loop invariant at line 25.

$$i = i \quad (26)$$

$$j = j \leq i \quad (27)$$

$$\text{value} = T_{m+1}(j, \gamma^j(m+1)) \quad (28)$$

By the loop invariant, we have the following after the loop at line 28:

$$i = i \quad (29)$$

$$j = j = i \quad (30)$$

$$\text{value} = T_{m+1}(i, \gamma^i(m+1)) \quad (31)$$

which concludes, since $i = l$.

Now we only need to prove the second loop invariant. First, at the beginning of the first iteration, we have $j = 0$ and **value** = $v = T_{m+1}(0, m+1)$ by (19), which satisfies the loop invariant. Now, assume that the invariant holds at the beginning of the j^{th} iteration. Here we need to consider only the case of $j < i$,

since the loop immediately terminates when $j = i$. Then, $j' = j + 1 \leq i$, and:

$$\begin{aligned}
T_{m+1}(j+1, \uparrow^{j+1}(m+1)) &= \text{hash}(T_{m+1}(j, \uparrow^j m), T_{m+1}(j, \uparrow^j(m+1))) \\
&\quad \text{(by Equation 10)} \\
&= \text{hash}(T_m(j, \uparrow^j m), \text{value}) \\
&\quad \text{(by Lemmas 1 \& 2 and Equation 28)} \\
&= \text{hash}(\text{branch}[j], \text{value}) \quad \text{(by Equations 16 \& 10)} \\
&= \text{value}'
\end{aligned}$$

Thus, the loop invariant holds at the beginning of the $(j+1)^{\text{th}}$ 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. $\text{zerohashes}[k] = Z(k)$
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 `zerohashes` 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 chain $\{T_m(k, \uparrow^k(m+1))\}_k$ 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 34, which suffices to conclude the main claim.

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

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 h' , size' , and root' denote the values at the next iteration $k = l+1$. Obviously, we have $h' = l+1$ and $\text{size}' = \uparrow^{l+1} m$. Also, we have $(\uparrow^l m) + 1 = \uparrow^l(m+1)$ by Lemma 1. Now, we have two cases:

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

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

- Case `size = 2l m` is even. Then, $2^l (m + 1)$ is odd. Thus,

$$\begin{aligned}
T_m(l + 1, 2^{l+1} (m + 1)) &= \text{hash}(T_m(l, 2^l (m + 1)), T_m(l, (2^l (m + 1)) + 1)) \\
&= \text{hash}(\text{root}, Z(l)) && \text{(by Lemma 3)} \\
&= \text{hash}(\text{root}, \text{zerohashes}[l]) && \text{(by Lemma 8)} \\
&= \text{root}'
\end{aligned}$$

Thus, we have $\text{root}' = T_m(l + 1, 2^{l+1} (m + 1))$, which concludes.

3.1 Refactoring Suggestion

The `deposit` function can be refactored as follows:

```

1  fun deposit'(value: int) -> unit:
2      assert deposit_count < 2^TREE_HEIGHT - 1
3      deposit_count += 1
4      size: int = deposit_count
5      i: int = 0
6      while i < TREE_HEIGHT - 1:
7          if size % 2 == 1:
8              break
9          value = hash(branch[i], value)
10         size /= 2
11         i += 1
12     branch[i] = value

```

The suggested refactoring is based on the following observations:

- Since `deposit` rejects when `deposit_count ≥ 2TREE_HEIGHT − 1`, the first loop of `deposit` cannot reach the last loop iteration, thus the loop bound can be safely decreased to `TREE_HEIGHT − 1`.
- The two loops of `deposit` can be combined into a single loop, which is slightly more efficient.

References

1. Ethereum Foundation: Deposit Contract in Eth 2.0. https://github.com/ethereum/deposit_contract/blob/master/deposit_contract/contracts/validator_registration.v.py
2. Ethereum Foundation: Ethereum 2.0 Specifications. <https://github.com/ethereum/eth2.0-specs>
3. 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. pp. 369–378. CRYPTO '87, Springer-Verlag, London, UK, UK (1988), <http://dl.acm.org/citation.cfm?id=646752.704751>
4. NIST: Perfect Binary Tree. <https://xlinux.nist.gov/dads/HTML/perfectBinaryTree.html>
5. Vitalik Buterin: Progressive Merkle Tree. https://github.com/ethereum/research/blob/master/beacon_chain_impl/progressive_merkle_tree.py