

# RED QUEEN’S SYNC PROTOCOL FOR ETHEREUM

ANDREW ASHIKHMIN & ALEXEY AKHUNOV

ABSTRACT. TODO: abstract.

”A slow sort of country!” said the Queen.  
”Now, here, you see, it takes all the running  
you can do, to keep in the same place. If you  
want to get somewhere else, you must run at  
least twice as fast as that!”

---

Lewis Carroll, Through the Looking-Glass  
and What Alice Found There

## 1. INTRODUCTION

In Red Queen’s Synchronisation Protocol for Ethereum 1x seeders reply with data as of their most recent block. That results in an inconsistent trie on the leecher initially (”phase 1”), which we patch later (”phase 2”). The idea is similar to that of Leaf Sync (Swende [2019]).

TODO: mention the sync failure problem Akhunov [2019a] and the needs of light clients like Mustekala. Inspirations like BitTorrent, Parity’s warp sync, Firehose Sync, Light Client Protocol.

N.B. Snapshot synchronisation rather than from genesis.

## 2. NOTATION

We mostly follow the conventions and notations of the Yellow Paper (Wood [2018]), for instance  $\mathbb{Y}$  denotes the set of nibble sequences. We use the letter  $\pi$  for prefixes of state or storage trie keys  $\mathbf{k} \in \mathbb{B}_{32}$ ,

$$(1) \quad \pi \in \mathbb{Y} \wedge \|\pi\| \leq 64$$

A key matches a prefix iff all their first nibbles are the same,

$$(2) \quad \text{MATCH}(\mathbf{k}, \pi) \equiv \forall_{i < \|\pi\|} : \mathbf{k}'[i] = \pi[i]$$

( $\mathbf{k}'$  is a sequence of nibbles, while  $\mathbf{k}$  is a sequence of bytes.)

## 3. PROTOCOL SPECIFICATION

We propose the following 3 request/reply operative pairs.

**GetStorageSizes** [+0x20, reqID:  $\mathbb{N}$ , blockAtLeast:  $\mathbb{N}$ , [account<sup>0</sup>:  $\mathbb{B}_{20}$ , account<sup>1</sup>:  $\mathbb{B}_{20}$ , ...]] Request storage trie sizes as of block #**blockAtLeast** or newer.

**StorageSizes** [+0x21, reqID:  $\mathbb{N}$ , blockNumber:  $\mathbb{N}$ , [numLeaves<sup>0</sup>:  $\mathbb{N}|\emptyset$ , numLeaves<sup>1</sup>:  $\mathbb{N}|\emptyset$ , ...]] Reply to **GetStorageSizes**. Returns storage trie sizes as of block #**blockNumber**  $\geq$  **blockAtLeast**. The elements returned must strictly match the accounts requested. The peer may return the empty list  $\emptyset$  instead of the number of leaves for accounts it does not have enough information about.

**GetNodeData2** [+0x22, reqID:  $\mathbb{N}$ , blockAtLeast:  $\mathbb{N}$ , [account<sup>0</sup>:  $\mathbb{B}_{20}|\emptyset$ , prefix<sup>0</sup>:  $\mathbb{Y}$ , prefix<sup>1</sup>:  $\mathbb{Y}$ , ...], [account<sup>1</sup>:  $\mathbb{B}_{20}|\emptyset$ , prefix<sup>0</sup>:  $\mathbb{Y}$ , prefix<sup>1</sup>:  $\mathbb{Y}$ , ...], ...] Request state or storage trie nodes as of block #**blockAtLeast** or newer. The empty list  $\emptyset$  instead of the account address signifies the state (rather than storage) trie. Note that this operative is similar to **GetNodeData** from Ethereum Wire Protocol PV63, but it uses prefixes rather than hashes as node keys<sup>1</sup>. TODO: prefix–node correspondence is trivial for branch nodes, not so much for leaf or extension nodes. TODO: prefix encoding consistent with the Yellow Paper.

---

*Date:* March 2019.

<sup>1</sup>For a justification see Péter Szilágyi’s comment at ETH v64 Wire Protocol Ring.

**NodeData2** [+0x23, reqID:  $\mathbb{N}$ , blockNumber:  $\mathbb{N}$ , [node<sub>0</sub><sup>0</sup>:  $\mathbb{B}$ , node<sub>1</sub><sup>0</sup>:  $\mathbb{B}$ , ...], [node<sub>0</sub><sup>1</sup>:  $\mathbb{B}$ , node<sub>1</sub><sup>1</sup>:  $\mathbb{B}$ , ...], ...] Reply to **GetNodeData2**. Returns trie nodes as of block #**blockNumber**  $\geq$  **blockAtLeast**. The nodes returned must strictly match the prefixes requested. The empty list  $\emptyset$  returned instead of a node means that the peer does not have enough information about the node requested.

**GetSubtries** [+0x24, reqID:  $\mathbb{N}$ , blockAtLeast:  $\mathbb{N}$ ,  
 [account<sup>0</sup>:  $\mathbb{B}_{20}|\emptyset$ ,  
 [prefix<sub>0</sub><sup>0</sup>:  $\mathbb{Y}$ , fromLevel<sub>0</sub><sup>0</sup>:  $\mathbb{N}$ , subtrieHash<sub>0</sub><sup>0</sup>:  $\mathbb{B}_{32}|\emptyset$ ],  
 [prefix<sub>1</sub><sup>0</sup>:  $\mathbb{Y}$ , fromLevel<sub>1</sub><sup>0</sup>:  $\mathbb{N}$ , subtrieHash<sub>1</sub><sup>0</sup>:  $\mathbb{B}_{32}|\emptyset$ ],  
 ...  
 ],  
 [account<sup>1</sup>:  $\mathbb{B}_{20}|\emptyset$ ,  
 [prefix<sub>0</sub><sup>1</sup>:  $\mathbb{Y}$ , fromLevel<sub>0</sub><sup>1</sup>:  $\mathbb{N}$ , subtrieHash<sub>0</sub><sup>1</sup>:  $\mathbb{B}_{32}|\emptyset$ ],  
 [prefix<sub>1</sub><sup>1</sup>:  $\mathbb{Y}$ , fromLevel<sub>1</sub><sup>1</sup>:  $\mathbb{N}$ , subtrieHash<sub>1</sub><sup>1</sup>:  $\mathbb{B}_{32}|\emptyset$ ],  
 ...  
 ],  
 ...  
 ]

] Request state or storage subtrie leaves along with proof nodes as of block #**blockAtLeast** or newer. The empty list  $\emptyset$  instead of the account address signifies state rather than storage trie. **fromLevel** specifies the number of upper nodes to be excluded from the proof in case the chain has not moved ahead (reply block is not newer). **subtrieHash** is a means to avoid resending leaves that have not changed.

**Subtries** [+0x25, reqID:  $\mathbb{N}$ , blockNumber:  $\mathbb{N}$ ,  
 [  
 [[node<sub>00</sub><sup>0</sup>:  $\mathbb{B}$ , node<sub>01</sub><sup>0</sup>:  $\mathbb{B}$ , ...], tooManyLeaves<sub>0</sub><sup>0</sup>, [key<sub>00</sub><sup>0</sup>:  $\mathbb{B}_{32}$ , val<sub>00</sub><sup>0</sup>:  $\mathbb{B}$ , key<sub>01</sub><sup>0</sup>:  $\mathbb{B}_{32}$ , val<sub>01</sub><sup>0</sup>:  $\mathbb{B}$ , ...]<sub>opt</sub>],  
 [[node<sub>10</sub><sup>0</sup>:  $\mathbb{B}$ , node<sub>11</sub><sup>0</sup>:  $\mathbb{B}$ , ...], tooManyLeaves<sub>1</sub><sup>0</sup>, [key<sub>10</sub><sup>0</sup>:  $\mathbb{B}_{32}$ , val<sub>10</sub><sup>0</sup>:  $\mathbb{B}$ , key<sub>11</sub><sup>0</sup>:  $\mathbb{B}_{32}$ , val<sub>11</sub><sup>0</sup>:  $\mathbb{B}$ , ...]<sub>opt</sub>],  
 ...  
 ],  
 [  
 [[node<sub>00</sub><sup>1</sup>:  $\mathbb{B}$ , node<sub>01</sub><sup>1</sup>:  $\mathbb{B}$ , ...], tooManyLeaves<sub>0</sub><sup>1</sup>, [key<sub>00</sub><sup>1</sup>:  $\mathbb{B}_{32}$ , val<sub>00</sub><sup>1</sup>:  $\mathbb{B}$ , key<sub>01</sub><sup>1</sup>:  $\mathbb{B}_{32}$ , val<sub>01</sub><sup>1</sup>:  $\mathbb{B}$ , ...]<sub>opt</sub>],  
 [[node<sub>10</sub><sup>1</sup>:  $\mathbb{B}$ , node<sub>11</sub><sup>1</sup>:  $\mathbb{B}$ , ...], tooManyLeaves<sub>1</sub><sup>1</sup>, [key<sub>10</sub><sup>1</sup>:  $\mathbb{B}_{32}$ , val<sub>10</sub><sup>1</sup>:  $\mathbb{B}$ , key<sub>11</sub><sup>1</sup>:  $\mathbb{B}_{32}$ , val<sub>11</sub><sup>1</sup>:  $\mathbb{B}$ , ...]<sub>opt</sub>],  
 ...  
 ],  
 ...  
 ]

] Reply to **GetSubtries**. Returns subtrie leaves with proofs as of block #**blockNumber**  $\geq$  **blockAtLeast**. The subtries returned must strictly match the prefixes requested. If the peer does not have information regarding a particular subtrie, it should return the empty list  $\emptyset$  (e.g. []) rather than [nodes, tooManyLeaves, leaves] for it. The nodes returned are the upper nodes of the trie down to the subtrie root, so that it is possible to verify that the leaves do belong to the Merkle Patricia trie in question. The first **fromLevel** upper nodes must be skipped if and only if **blockNumber** = **blockAtLeast**. (If **fromLevel** = 0, then the nodes must start with the root node.) **tooManyLeaves** is a boolean flag (0 = *false*, 1 = *true*) indicating that the subtrie requested contains too many leaves. TODO: how many is too many? The leaves are represented as the list of their keys<sup>2</sup> and values. The peer may only return either all leaves of the subtrie or nothing. In case of **tooManyLeaves** the leaves should not returned<sup>3</sup>.

Proof nodes in **Subtrie** replies give us a means to detect faulty or malicious peers. Note that state trie replies do not inline storage tries unlike Leaf Sync.

TODO: is block number OK given chain reorgs?

#### 4. SUGGESTED FULL SYNC ALGORITHM

Here we suggest a possible algorithm for full state and storage snapshot synchronisation using the protocol specified above. We describe a modus operandi where the seeder replies with its most recent data, and the

<sup>2</sup>It is feasible to return suffixes rather than full keys given that prefixes are known, but we deem the performance gain to be insignificant.

<sup>3</sup>In that case the peer must not return an empty list as that would imply that there are no leaves matching the given prefix.

leecher has to handle trie data coming from different blocks. We suggest to perform synchronisation in two stages: during phase 1 the leecher obtains leaf data (with the necessary proof nodes) as of any reasonable block height, while during phase 2 it patches up the trie in order to catch up to the most recent block<sup>4</sup>. The idea was proposed in Swende [2019].

Let us focus on the state trie for the moment; we shall come back to storage sync later. For phase 1 we suggest to send **GetSubtries** requests with a single prefix per request, ditto for phase 2. All requested prefixes are of size  $d_1$  during phase 1 and of  $d_2$  during phase 2,  $d_2 \geq d_1$ . We elaborate on the values of  $d_1$  and  $d_2$  later. The leecher gradually builds the first upper  $d_2$  levels of the Merkle Patricia trie<sup>5</sup>. (The full trie can be constructed if so desired, but only the upper  $d_2$  levels are necessary for our algorithm.) Populated nodes are marked with the block number as of they are valid. The algorithm preserves the following invariant: parent's block is always no older than child's block.

During phase 1 the leecher requests each possible prefix of size  $d_1$  exactly once (barring network failures and faulty peers). When sending a request, the leecher sets its **blockAtLeast** to the block of the root of the current (partially populated) trie, **fromLevel** to the number of populated nodes down the path/prefix that are of the same block as the root, and **subtrieHash** is not used (set to  $\emptyset$ ). Having received a reply, the leecher verifies its proof. If the proof is valid, the leecher writes received leafs to the database and updates the nodes along the prefix/path. By the end of phase 1 the leecher will have all state accounts populated, albeit inconsistently.

TODO: an example.

TODO: phase 2.

TODO: storage tries.

## 5. PERFORMANCE ANALYSIS

For this analysis we assume that all tries are well balanced. We also assume that all top nodes up to a certain trie level  $i$  are branch nodes, not leaf nor extension nodes. This is a reasonable assumption if  $i$  is not too big—see Akhunov [2019b].

Some notation (TODO: internal notation consistency + cross-check against the Yellow Paper):

$o$  – **Leaves** reply overhead in bytes.

$b$  – size of a branch node in bytes.

$l$  – average leaf size in bytes.

$||R_b||$  – number of nodes in reply.

$||R_l||$  – number of leaves in reply.

Thus the size of a reply, assuming the average leaf size, is

$$(3) \quad S(R) = o + ||R_b||b + ||R_l||l$$

TODO: RLP changes the formula slightly.

$t$  – total number of leaves in the trie.

TODO: optimal phase 1 depth.

Now let us find the maximum size  $d$  of the request prefix  $\pi$  that makes sense to use when we are catching up (phase 2 of the sync). Let assume that only one leaf that matches  $\pi$  has changed. (If we know that there are no changes, there is no need for a sync request.) That is a reasonable assumption if we are not too many blocks behind, there are not that many changes per block, and  $d$  is sufficiently large. (For instance, if we are 100 blocks behind, and there are 500 leaf changes per block, then  $d \geq 4$  will suffice.) Consider two options: request the prefix  $\pi$  or send requests with prefixes  $\pi \cdot 0, \dots, \pi \cdot 15$  of size  $d + 1$  (not necessarily all 16 of them). In the first case we receive a reply of the size, on average,

$$(4) \quad S = o + ||R_b||b + \frac{t}{16^d}l$$

For the second case we need to send at most two requests, as the first reply will give us the information to identify which nibble has changed. (With the probability  $\frac{1}{16}$  the second request is not necessary.) The

<sup>4</sup>The Red Queen's race is a nice metaphor for phase 2.

<sup>5</sup> $d_2$  is small enough so that we can reasonably assume that (almost) all nodes in question are branch nodes; see Akhunov [2019b].

combined size of those 1 or 2 replies, on average, is

$$(5) \quad S' = \left(1 + \frac{15}{16}\right) o + (||R_b|| + 1)b + \frac{t}{16^{d+1}}l$$

It does not make sense to prefer requests with longer prefixes if  $S \leq S'$ . Solving this inequality, we obtain

$$(6) \quad 16^d(16b + 15o) \geq 15tl$$

In other words, it does not make sense to use prefixes longer than

$$(7) \quad \left\lceil \log_{16} \frac{15tl}{16b + 15o} \right\rceil$$

TODO: convergence analysis.

## 6. CONCLUSION

TODO: conclusion.

## REFERENCES

- Alexey Akhunov. Looking back at the Ethereum 1x workshop 26–28.01.2019 (part 1), January 2019a. URL <https://medium.com/@akhounov/looking-back-at-the-ethereum-1x-workshop-26-28-01-2019-part-1-70c1ebd93266>.
- Alexey Akhunov. Looking back at the Ethereum 1x workshop 26–28.01.2019 (part 2), February 2019b. URL <https://medium.com/@akhounov/looking-back-at-the-ethereum-1x-workshop-26-28-01-2019-part-2-d3d8fdcede10>.
- Martin Holst Swende. Leaf sync, March 2019. URL [https://notes.ethereum.org/kphcc\\_CKT4a5sUs\\_zWVela](https://notes.ethereum.org/kphcc_CKT4a5sUs_zWVela).
- Gavin Wood. Ethereum: A Secure Decentralised Generalised Transaction Ledger, December 2018. URL <https://github.com/ethereum/yellowpaper>.